



# Null space gradient flows for constrained optimization with applications to shape optimization

Florian Feppon, Grégoire Allaire, Charles Dapogny

## ► To cite this version:

Florian Feppon, Grégoire Allaire, Charles Dapogny. Null space gradient flows for constrained optimization with applications to shape optimization. 2019. hal-01972915v1

**HAL Id: hal-01972915**

**<https://hal.science/hal-01972915v1>**

Preprint submitted on 8 Jan 2019 (v1), last revised 16 Nov 2020 (v4)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# NULL SPACE GRADIENT FLOWS FOR CONSTRAINED OPTIMIZATION WITH APPLICATIONS TO SHAPE OPTIMIZATION

F. FEPPON<sup>12</sup>, G. ALLAIRE<sup>1</sup>, C. DAPOGNY<sup>3</sup>

<sup>1</sup> *Centre de Mathématiques Appliquées, École Polytechnique, Palaiseau, France*

<sup>2</sup> *Safran Tech, Magny-les-Hameaux, France*

<sup>3</sup> *Univ. Grenoble Alpes, CNRS, Grenoble INP<sup>1</sup>, LJK, 38000 Grenoble, France*

**ABSTRACT.** The purpose of this article is to introduce a gradient-flow algorithm for solving generic equality or inequality constrained optimization problems, which is suited for shape optimization applications. We rely on a variant of the Ordinary Differential Equation (ODE) approach proposed by Yamashita [48] for equality constrained problems: the search direction is a combination of a null space step and a range space step, which are aimed to reduce the value of the minimized objective function and the violation of the constraints, respectively. Our first contribution is to propose an extension of this ODE approach to optimization problems featuring both equality and inequality constraints. In the literature, a common practice consists in reducing inequality constraints to equality constraints by the introduction of additional slack variables. Here, we rather solve their local combinatorial character by computing the projection of the gradient of the objective function onto the cone of feasible directions. This is achieved by solving a dual quadratic programming subproblem whose size equals the number of active or violated constraints, and which allows to identify the inequality constraints which should remain tangent to the optimization trajectory. Our second contribution is a formulation of our gradient flow in the context of—infinite-dimensional—Hilbert space settings. This allows to extend the method to quite general optimization sets equipped with a suitable manifold structure, and notably to sets of shapes as it occurs in shape optimization with the framework of Hadamard’s boundary variation method. The cornerstone of this latter setting is the classical operation of extension and regularization of shape derivatives. Some numerical comparisons on simple academic examples are performed to illustrate the behavior of our algorithm. Its numerical efficiency and ease of implementation are finally demonstrated on more realistic shape optimization problems.

**Keywords.** nonlinear constrained optimization, gradient flows, shape and topology optimization, null space method.

**AMS Subject classifications.** 65K10, 49Q10, 34C35, 49B36, 65L05.

---

## CONTENTS

1. Introduction	2
2. Gradient flows for equality-constrained optimization in Hilbert spaces	5
2.1. Notations and first-order optimality conditions	5
2.2. Definitions and properties of the null space and range space steps $\xi_J$ and $\xi_C$	6
2.3. Behavior of the trajectories of the flow	9
3. Proposed extension to equality and inequality constraints	11
3.1. Notations and preliminaries	11
3.2. The method of slack variables for inequality constraints	12
3.3. The proposed algorithm	13
3.4. Comparison between the proposed method and the use of slack variables	20
4. Practical implementation details	21
4.1. Time step adaptation based on a merit function.	21
4.2. Accounting for discontinuities near the inequality constraint barriers	22
5. Illustrations and comparisons on academic test cases	23

---

Corresponding author. Email: [florian.feppon@polytechnique.edu](mailto:florian.feppon@polytechnique.edu).

<sup>1</sup>Institute of Engineering Univ. Grenoble Alpes

5.1. Test case 1 : unfeasible initialization with initial gradient aligned with the constraints.	24
5.2. Test case 2 : unfeasible initialization with initial gradient not aligned with the constraints.	25
5.3. Test case 3: a saturated inequality constraint becoming inactive along the optimization path	27
<b>6. Optimization on smooth manifolds: application to shape optimization</b>	<b>29</b>
6.1. Hadamard’s framework for gradient based shape optimization	30
6.2. Implementation of the constrained gradient flow for level set based shape optimization	32
6.3. Illustrations on a multiple load structural shape optimization test case	33
<b>7. Conclusion and perspectives</b>	<b>38</b>
<b>References</b>	<b>41</b>

---

## 1. INTRODUCTION

Over the past decades, many iterative algorithms have been designed for generic constrained optimization problems of the form:

$$\begin{aligned} \min_{x \in \mathcal{X}} \quad & J(x) \\ \text{s.t.} \quad & \begin{cases} \mathbf{g}(x) = 0 \\ \mathbf{h}(x) \leq 0, \end{cases} \end{aligned} \quad (1.1)$$

where  $\mathcal{X}$  is the optimization set,  $J : \mathcal{X} \rightarrow \mathbb{R}$  is a differentiable objective function,  $\mathbf{g} : \mathcal{X} \rightarrow \mathbb{R}^p$  and  $\mathbf{h} : \mathcal{X} \rightarrow \mathbb{R}^q$  are differentiable functions accounting for  $p$  equality and  $q$  inequality differentiable constraints, respectively. ‘Classical’ gradient-based algorithms for the numerical resolution of (1.1) include, e.g., Penalty, Lagrangian, Interior Point and Trust Region Methods, Sequential Quadratic or Linear Programming (SQP or SLP) [14, 36, 49, 29], the Method of Moving Asymptotes (MMA) [44], the Method of Feasible Directions [51]. When dealing with a given application, the choice of one particular optimization method is partly determined by the difficulty of implementation. As is often the case in constrained optimization, it turns out that all the aforementioned techniques require fine tuning of the algorithm parameters in order to achieve convergence with a reasonable computational efficiency; these parameters are e.g. the penalty coefficients in the Augmented Lagrangian and Interior Point methods, the size of the trust region in SLP algorithms, the strategy for approximating the Hessian matrix in SQP, the bounds on the asymptotes in MMA and the Topkis parameters in MFD. The correct determination of these parameters is strongly case-dependent and often unintuitive: for instance, the penalty coefficients must be neither ‘too large’ nor ‘too small’ in Lagrangian methods, the SLP trust region size—which acts as a step length—cannot be chosen too small (otherwise the involved quadratic subproblems may not have a solution).

In a slightly different spirit, a significant amount of work has been devoted to dynamical system approaches for addressing problems of the form (1.1): a solution is reached as the stationary point  $x^*$  of a continuous trajectory  $x(t)$ , solving a suitable Ordinary Differential Equation (ODE). When  $\mathcal{X} = \mathbb{R}^k$  is a finite-dimensional vector space and in the particular, unconstrained case where constraints are omitted in (1.1), the most basic of these is the celebrated gradient flow:

$$\dot{x}(t) = -\nabla J(x(t)). \quad (1.2)$$

The value of the objective  $t \mapsto J(x(t))$  is guaranteed to decrease along the optimization path, and the Euler step size  $\Delta t$ , which is the unique parameter involved in the discretization of (1.2), can always be made sufficiently small so that a reduction of the value of the objective function is observed. In this sense, such dynamical system approaches can be expected to be more reliable than the aforementioned iterative algorithms, although they might not achieve the fastest rate of convergence. Extensions of the standard gradient flow (1.2) have been proposed to account for equality constraints  $\mathbf{g}(x) = 0$ . For instance, Tanabe [45] proposed to replace the gradient  $\nabla J(x)$  in (1.2) by its tangential projection  $\boldsymbol{\xi}_J(x)$  onto the feasible set. Yamashita [48] suggested to add a Gauss-Newton direction  $\boldsymbol{\xi}_C(x)$  to this ODE in order to smoothly lead the optimization path toward the feasible domain; the resulting dynamical system for equality constrained optimization reads:

$$\dot{x}(t) = -\alpha_J \boldsymbol{\xi}_J(x(t)) - \alpha_C \boldsymbol{\xi}_C(x(t)), \quad (1.3)$$

where  $\xi_J$  and  $\xi_C$  are respectively defined by

$$\xi_J(x) = (I - Dg^T(DgDg^T)^{-1}Dg)\nabla J(x), \quad (1.4)$$

$$\xi_C(x) = Dg^T(DgDg^T)^{-1}g(x), \quad (1.5)$$

with  $I$  the identity mapping and  $(Dg(x))_{ij} = \partial_j g_i(x)$  the Jacobian matrix of the constraint function  $g(x)$  (we shall omit the dependence with respect to  $x$  when the context is clear). The descent direction  $\dot{x}$  is a combination of the so-called ‘null space’ direction  $\xi_J(x) \in \text{Ker}(Dg(x))$  and ‘range space’ direction  $\xi_C(x) \in \text{span}(Dg^T(x))$ , lying respectively in the null space of the constraints and in its orthogonal complement. For this reason, we call the ODE (1.3) a ‘null space’ gradient flow. In (1.3),  $\alpha_J, \alpha_C > 0$  are two (facultative) parameters that we choose to introduce in order to scale how fast the objective function should be decreased which respect to the violation of the constraints. This decomposition of the descent direction ensures that the constraints vanish at an exponential rate, namely  $g(x(t)) = e^{-\alpha_C t}g(x_0)$ , while the objective function decreases along the direction  $-\xi_J(x(t))$  (see [48] and Proposition 1 hereafter). The discretization of the flow (1.3) can be related to SQP (see [38]) and to null space iterative methods [13, 14, 36] (which, to our knowledge, apply only to equality-constrained optimization).

When it comes to handling inequality constraints, the most common approach consists in introducing  $q$  slack variables  $\{z_i\}_{1 \leq i \leq q} \in \mathbb{R}^q$  so as to convert the  $q$  inequalities  $h_i(x) \leq 0$  for  $1 \leq i \leq q$  into as many equality constraints  $h_i(x) + z_i^2 = 0$ , before then solving the ODE (1.3) in the augmented space  $(x, z) \in \mathbb{R}^k \times \mathbb{R}^q$ . It is even possible to eliminate the variables  $z_i$  from the latter ODE if the initialization is feasible, see [32, 41]. Schropp and Singer showed in [38] that stationary points of the resulting dynamical system fulfill partially the Karush Kuhn Tucker (KKT), necessary first-order optimality condition: at such a stationary point  $x^*$ , there exist multipliers  $(\lambda(x^*), \mu(x^*)) \in \mathbb{R}^p \times \mathbb{R}^q$  such that

$$\nabla J(x^*) + Dg(x^*)^T \lambda(x^*) + Dh(x^*)^T \mu(x^*) = 0, \text{ and } h_i(x^*)\mu_i(x^*) = 0 \text{ for } i = 1, \dots, q.$$

However, the complete set of KKT conditions may not be satisfied because of possible negative values of some components of the Lagrange multiplier  $\mu_i(x^*) < 0$ . Nevertheless, this strategy proves efficient in practice because among all possible stationary points, only those satisfying the complete KKT conditions are asymptotically stable; see again [38] about this point.

The contributions of the present article are twofold. First, we propose an extension of the flow (1.3) which is able to account for inequality constraints without the need to increase the size of the problem with auxiliary slack variables. Our approach relies on a suitable dual program for the resolution of the combinatorial character of the inequality constraints: we precisely identify to which subset of the active inequality constraints the optimization path is allowed to ‘unstick’ (thus re-entering into the feasible domain) and, conversely, to which inequality constraints it must remain tangent. The resulting dynamical system is slightly more robust than those based on slack variables, since we shall see that all its stationary points are true KKT points of (1.1), and conversely. More specifically, for a given subset of indices  $I \subset \{1, \dots, q\}$ , denote  $\mathbf{h}_I(x) := (h_i(x))_{i \in I}$  the corresponding inequality constraints and  $\mathbf{C}_I(x)$  the matrix

$$\mathbf{C}_I(x) := \begin{bmatrix} g(x) \\ \mathbf{h}_I(x) \end{bmatrix}. \quad (1.6)$$

Then, for inequality constrained problems, we define new directions  $\xi_J(x)$  and  $\xi_C(x)$  in (1.3) as follows:

$$\xi_J(x) = (I - DC_{\tilde{I}(x)}^T(DC_{\tilde{I}(x)}DC_{\tilde{I}(x)}^T)^{-1}DC_{\tilde{I}(x)})\nabla J(x), \quad (1.7)$$

$$\xi_C(x) = DC_{\tilde{I}(x)}^T(DC_{\tilde{I}(x)}DC_{\tilde{I}(x)}^T)^{-1}\mathbf{C}_{\tilde{I}(x)}(x), \quad (1.8)$$

which involve two different subsets  $\hat{I}(x) \subset \tilde{I}(x) \subset \{1, \dots, q\}$ :  $\tilde{I}(x)$  is the set of all saturated or violated constraints, defined by

$$\tilde{I}(x) = \{i \in \{1, \dots, q\} \mid h_i(x) \geq 0\}, \quad (1.9)$$

and  $\hat{I}(x)$  is a subset of  $\tilde{I}(x)$  identifying the constraints onto which the gradient of the objective function is projected tangentially. The set  $\hat{I}(x)$  shall be obtained by solving a dual quadratic optimization subproblem (equation (3.13) below) of the size of  $\tilde{I}(x)$ . As detailed later in Proposition 4, this dual problem ensures that  $\xi_J(x)$  is the projection of the gradient  $\nabla J(x)$  onto the cone of feasible directions. As a result,  $-\xi_J(x)$  is

always the best possible descent direction respecting locally both equality and inequality constraints. Since the sets  $\tilde{I}(x)$  or  $\hat{I}(x)$  are subject to change as soon as inequality constraints become active or inactive, or if not enough regularity holds, the ODE (1.3) has a discontinuous right-hand side and is defined only formally (note that a rigorous mathematical meaning could still be provided with the theory of non smooth ODEs, see [21, 28]). However and as we shall detail further on, its discretization makes sense and exhibits the same decreasing properties as its continuous counterpart for sufficiently small steps  $\Delta t$ .

Our second main contribution is the exposure of our dynamical system strategy in a setting that can handle quite general ranges of infinite-dimensional optimization sets, including for instance those involved in non parametric shape and topology optimization, which is our final goal. Here, a clear distinction needs to be made between the Fréchet derivative and the gradient of the objective and constraint functionals: if  $\mathcal{X} = V$  is a Hilbert space, recall that the gradient is obtained by identifying the differential (a linear form on  $V$ ) to an element in  $V$  via the Riesz representation theorem; see Definition 1. In order to account for the infinite dimensionality, we shall specify clearly how (1.7) and (1.8) must be computed with respect to such identification and used when the optimization set  $\mathcal{X}$  is a Hilbert space or a more general set equipped with a suitable manifold structure. As we have mentioned, our ultimate motivations actually originate from the field of shape optimization based on the method of Hadamard, for which the minimization set  $\mathcal{X}$  in (1.1) is the set of all possible open Lipschitz subdomains  $\Omega$  enclosed in some ‘hold-all’ domain  $D \subset \mathbb{R}^d$ :

$$\mathcal{X} = \{\Omega \subset D \mid \Omega \text{ Lipschitz}\}. \quad (1.10)$$

This set is not a vector space, but it can be locally parameterized by the Sobolev space  $W^{1,\infty}(D, \mathbb{R}^d)$ . In order to minimize a shape functional  $\Omega \mapsto J(\Omega)$ , we determine the best local variation of the form  $(I + \theta)(\Omega)$  where  $\theta \in W^{1,\infty}(D, \mathbb{R}^d)$  can be interpreted as a (sufficiently small) displacement field [35, 42, 30]. As we shall review later on in Section 6.1, this endows  $\mathcal{X}$  with a manifold structure which allows for gradient based optimization [6, 47]. In shape optimization, the identification of the gradient is achieved by solving an extension and regularization problem, which has some very important consequences in numerical algorithms, see e.g. [17, 19]. We shall detail below how this step is naturally and consistently included in our algorithm, a matter which so far does not seem fully clear in the literature concerned with constrained shape optimization: common approaches rather compute a descent direction *first*, before performing a regularization, see e.g. [23, 25].

In the open academic literature on shape optimization based on Hadamard’s method, advanced mathematical programming methods are not frequently described. Rather, for simplicity of implementation, Penalty and Augmented Lagrangian Methods are often used, all the more when only one constraint is considered [6, 18]. Morin et. al. introduced a variant of SQP in [34] but treated a volume constraint with a Lagrange Multiplier method. For more complex applications, some authors have introduced adapted variants of Sequential Linear Programming [23] or of the Method of Feasible Direction [25]. However, the high dimensionality and complexity of the research space leaves very little intuition on how to select appropriately the parameter values featured in the implementation of these latter algorithms. As a result, a fair amount of trials and errors is often required in order to obtain satisfying minimizing sequences of shapes—a process which can be very time consuming, especially when every optimization step depends on the resolution of some physical model involving partial differential equations. As a result of the flexibility of the ODE approaches, our method depends truly only on the discretization step  $\Delta t$ , and on the physically interpretable dimensionless parameters  $\alpha_J, \alpha_C$ , which makes them relatively easy to tune for the user.

Several contributions in the field of shape and topology optimization can be related to ours. In fact, our method is very close in spirit to the recent work of Barbarosie et. al. [12], who derived an iterative algorithm for equality constrained optimization which turns out to be a discretization of (1.3) with a variable scaling for the parameter  $\alpha_C$ . For inequality constraints, the authors proposed (without convergence results) an active set strategy also based on the extraction of an appropriate subset of the active constraints. However their method relies on a different algorithm from ours, that yields generally a different set than  $\hat{I}(x)$ , see Remark 4 below for more details. Finally, Yulin and Xiaoming [50] also suggested to project the gradient of the objective function onto the convex cone of feasible directions; nevertheless, they remained elusive regarding how the projection problem is solved or how violated constraints are tackled.

The present article is organized as follows. In [Section 2](#), we review the definition and the properties of the gradient flow [\(1.3\)](#) for equality constrained optimization in the case where the minimisation set  $\mathcal{X}$  is a Hilbert space. We detail then in [Section 3](#) the necessary adaptations to account for inequality constraints and in particular the introduction of the dual subproblem allowing to determine the null space direction  $\xi_J(x)$ . Under some technical assumptions, we prove in [Proposition 5](#) the convergence of the “null space” gradient flow [\(1.3\)](#) towards points satisfying the full KKT condition. We detail algorithmic implementation aspects in [Section 4](#). In [Section 5](#), we provide pedagogical illustrations of our method on simple academic test cases, and we compare it to the method of slack variables for inequality constraints. We finally focus on shape optimization applications in [Section 6](#). After clarifying the necessary adaptations required to extend the discretization of [\(1.3\)](#) to the set  $\mathcal{X}$  defined by [\(1.10\)](#), we explain how our algorithm can be integrated within the level set method for shape optimization [\[47, 6, 5\]](#). We then implement and demonstrate numerically this method on the optimal design of a bridge structure subject to multiple loads, which involves up to ten constraints. A conclusion and several perspectives are outlined in [Section 7](#).

## 2. GRADIENT FLOWS FOR EQUALITY-CONSTRAINED OPTIMIZATION IN HILBERT SPACES

In this section, we consider the case where the optimization takes place on a Hilbert space  $\mathcal{X} = V$  with inner product  $a(\cdot, \cdot)$  and relative norm  $\|\cdot\|_V = a(\cdot, \cdot)^{1/2}$ ; see [Section 6](#) for the description of the more general situation associated to our shape optimization applications. The first focus of our study is the minimization problem [\(1.1\)](#) where only equality constraints are present, namely:

$$\begin{aligned} \min_{x \in V} \quad & J(x) \\ \text{s.t.} \quad & \mathbf{g}(x) = 0, \end{aligned} \tag{2.1}$$

where  $J : V \rightarrow \mathbb{R}$  and  $\mathbf{g} : V \rightarrow \mathbb{R}^p$  are Fréchet differentiable functions. Our purpose is to recall how the ODE approach [\(1.3\)](#) can be extended to the present Hilbertian setting. Let us emphasize that, although this section is elementary and not completely new, it is not easily found as is in the literature. Since it is key in understanding our technique for handling inequality constraints in [Section 3](#), the present context is thoroughly detailed for the reader’s convenience.

The section is organized as follows. We first recall the definitions of the differential, the gradient, and the transpose of the differential in the Hilbertian context. We then sketch briefly how the formulas [\(1.4\)](#) and [\(1.5\)](#) can be formally obtained. We state the properties of the null space step  $\xi_J(x)$  and its relation to Lagrange multipliers by means of a dual problem in [Lemma 1](#). Finally, we review the decrease properties of the obtained dynamical system in [Proposition 1](#).

### 2.1. Notations and first-order optimality conditions

We start by setting notations about differentiability and gradients in Hilbert spaces that we use throughout this article. Our notations may differ from those used by other authors because we need a clear distinction between gradient and Fréchet derivatives.

#### Definition 1.

- (1) A vector-valued function  $\mathbf{g} : V \rightarrow \mathbb{R}^p$  is *differentiable* at a point  $x \in V$  if there exists a continuous linear mapping  $D\mathbf{g}(x) : V \rightarrow \mathbb{R}^p$  such that

$$\mathbf{g}(x+h) = \mathbf{g}(x) + D\mathbf{g}(x)h + o(h) \text{ with } \frac{o(h)}{\|h\|_V} \xrightarrow{h \rightarrow 0} 0. \tag{2.2}$$

$D\mathbf{g}(x)$  is called the Fréchet derivative of  $\mathbf{g}$  at  $x$ .

- (2) If  $\mathbf{g} : V \rightarrow \mathbb{R}^p$  is differentiable, for any  $\mu \in \mathbb{R}^p$ , the Riesz representation theorem [\[15\]](#) ensures the existence of a unique vector  $D\mathbf{g}(x)^T \mu \in V$  satisfying

$$\forall \mu \in \mathbb{R}^p, \forall \xi \in V, a(D\mathbf{g}(x)^T \mu, \xi) = \mu^T D\mathbf{g}(x) \xi, \tag{2.3}$$

where the superscript  $T$  stands for the usual transpose of a vector in the Euclidean space  $\mathbb{R}^p$ . The linear operator  $D\mathbf{g}(x)^T : \mathbb{R}^p \rightarrow V$  thus defined is called the *transpose* of  $D\mathbf{g}(x)$ .

- (3) If  $J : V \rightarrow \mathbb{R}$  is a scalar function differentiable at  $x \in V$ , the Riesz representation theorem ensures the existence of a unique vector  $\nabla J(x) \in V$  satisfying

$$\forall \xi \in V, a(\nabla J(x), \xi) = DJ(x)\xi. \quad (2.4)$$

This vector  $\nabla J(x)$  is called the *gradient* of  $J$  at  $x$ .

Throughout the paper, we shall sometimes omit the explicit mention to  $x$  in the notations for differentials or gradients when the considered point  $x \in V$  is clear, so as to keep expressions as light as possible,

*Remark 1.*

- (1) If  $V$  is the (finite-dimensional) Euclidean space  $\mathbb{R}^k$ , equipped with the standard inner product, the Fréchet derivative and the transpose of a vector function  $\mathbf{g} : \mathbb{R}^k \rightarrow \mathbb{R}^p$  are respectively given by the Jacobian matrix  $(D\mathbf{g})_{ij} = \partial_j g_i$  and its transpose  $(D\mathbf{g}^T)_{ij} = (D\mathbf{g})_{ji} = \partial_i g_j$ . In the literature, the differential matrix  $D\mathbf{g}$  is often denoted with the nabla notation  $\nabla \mathbf{g}$ . For the sake of clarity, we reserve the  $\nabla$  symbol to denote the gradient of scalar functions  $J : V \rightarrow \mathbb{R}$ , and it holds that  $\nabla J(x) = DJ(x)^T \mathbf{1}$ . The curly transpose notation  $\mathcal{T}$  appearing in the objects  $DJ(x)^{\mathcal{T}}$  or  $D\mathbf{g}(x)^{\mathcal{T}}$  encodes at the same time the operator transposition (reversing the input and range spaces) and the Riesz identifications.
- (2) Still in the case where  $V = \mathbb{R}^k$  is finite-dimensional and  $a$  is given by a symmetric definite positive matrix  $A$  (that is  $a(\xi, \xi) = \xi^T A \xi$ ), the transpose of a  $p \times k$  matrix  $M : \mathbb{R}^k \rightarrow \mathbb{R}^p$  with respect to  $a$  is  $M^{\mathcal{T}} = A^{-1} M^T$ . As we shall see in [Section 6](#), in shape optimization applications,  $a$  often stands for the bilinear form associated to an elliptic operator, hence the calligraphic transpose  $\mathcal{T}$  encompasses the extension and regularization step of the shape derivative, see [Section 6.1](#) below. If  $V$  is replaced by the tangent space to some Riemannian manifold, the bilinear form  $a$  can be interpreted as a metric and  $\nabla J(x)$ , as given by (2.4), is the covariant gradient with respect to this metric.
- (3) When  $V$  is a general Hilbert space, for a vector-valued function  $\mathbf{g} : V \rightarrow \mathbb{R}^p$  with coordinates  $\mathbf{g}(x) = (g_i(x))_{1 \leq i \leq p}$ ,  $D\mathbf{g} : V \rightarrow \mathbb{R}^p$  is the ‘row’ matrix whose entries are the  $p$  linear forms  $Dg_i(x) : V \rightarrow \mathbb{R}$ . The transpose  $D\mathbf{g}(x)^{\mathcal{T}}$  is the ‘column’ matrix gathering the  $p$  vectors  $(\nabla g_i(x))_{1 \leq i \leq p}$  obtained by solving the  $p$  identification problems:

$$\forall \xi \in V, a(\nabla g_i(x), \xi) = Dg_i(x)\xi; \quad (2.5)$$

more precisely:

$$\forall \mu \in \mathbb{R}^d, D\mathbf{g}(x)^{\mathcal{T}} \mu = \sum_{i=1}^p \mu_i \nabla g_i(x).$$

In particular, the  $p \times p$  matrix  $D\mathbf{g}D\mathbf{g}^{\mathcal{T}} \in \mathbb{R}^{p \times p}$  has entries

$$(D\mathbf{g}D\mathbf{g}^{\mathcal{T}})_{ij} = a(\nabla g_i, \nabla g_j) = Dg_i(x)(\nabla g_j(x)).$$

Throughout this section, the equality constraints are said to be qualified at a point  $x \in V$  if

$$\text{rank}(D\mathbf{g}(x)) = p, \text{ or equivalently } D\mathbf{g}(x)D\mathbf{g}(x)^{\mathcal{T}} \text{ is an invertible } p\text{-by-}p \text{ matrix.} \quad (2.6)$$

Note that (2.6) makes sense even at points  $x \in V$  where  $\mathbf{g}(x) \neq 0$ , a fact that we shall use in the sequel. Under the above notations, let us recall the classical first-order necessary optimality conditions (KKT) for the equality-constrained problem (2.1) at some point  $x^* \in V$  [14, 36] where the constraints are satisfied and qualified: there exists  $\lambda(x^*) \in \mathbb{R}^p$  such that,

$$\begin{cases} \nabla J(x^*) + D\mathbf{g}(x^*)^{\mathcal{T}} \lambda(x^*) = 0, \\ \mathbf{g}(x^*) = 0. \end{cases} \quad (2.7)$$

## 2.2. Definitions and properties of the null space and range space steps $\xi_J$ and $\xi_C$

The defining formulas (1.4) and (1.5) for the steps  $\xi_J(x)$  and  $\xi_C(x)$ , featured in the ODE (1.3), can be given a meaning in the present Hilbert space setting owing to the definition (2.3) of the transpose:



**Definition 2.** Consider the optimization problem (2.1). For any point  $x \in V$  satisfying the constraint qualification condition (2.6), we define the *null space* and *range space* directions  $\xi_J(x)$  and  $\xi_C(x)$  by, respectively:

$$\xi_J(x) := (\mathbf{I} - \mathbf{Dg}^T(\mathbf{DgDg}^T)^{-1}\mathbf{Dg})\nabla J(x), \quad (2.8)$$

$$\xi_C(x) := \mathbf{Dg}^T(\mathbf{DgDg}^T)^{-1}\mathbf{g}(x). \quad (2.9)$$

Let us recall a formal intuition motivating the expressions (2.8), (2.9) and the dynamical system (1.3). Following the derivation of e.g. [11] and as is classical in Lagrange multiplier methods for optimization judging from the KKT optimality conditions (2.7), it is natural to search for an iterative optimization scheme (indexed by the subscript  $n$ ) of the form

$$x_{n+1} = x_n - \Delta t(\alpha_J \nabla J(x_n) + \mathbf{Dg}(x_n)^T \lambda_n), \quad (2.10)$$

where  $\lambda_n \in \mathbb{R}^p$  is a tentative value for the Lagrange multiplier  $\lambda$  in (2.7),  $\alpha_J$  is a user-defined coefficient and  $\Delta t$  is the step increment between successive iterations. We determine the value of  $\lambda_n$  by imposing that the constraint  $\mathbf{g}(x_{n+1})$  decreases by a factor  $1 - \alpha_C \Delta t$  at the next iteration, up to some first-order error in  $\Delta t$ . Since

$$\mathbf{g}(x_{n+1}) = \mathbf{g}(x_n) - \Delta t \mathbf{Dg}(x_n)(\alpha_J \nabla J(x_n) + \mathbf{Dg}(x_n)^T \lambda_n) + o(\Delta t),$$

the requirement that  $\mathbf{g}(x_{n+1}) \simeq (1 - \alpha_C \Delta t)\mathbf{g}(x_n)$  suggests the rule:

$$\lambda_n = (\mathbf{Dg}(x_n)\mathbf{Dg}(x_n)^T)^{-1}(\alpha_C \mathbf{g}(x_n) - \alpha_J \mathbf{Dg}(x_n)\nabla J(x_n)). \quad (2.11)$$

We recognize then the scheme (2.25) (a time discretization of (1.3)) by replacing  $\lambda_n$  with the above value (2.11) in (2.10).

### 2.2.1. Properties of the null space step $\xi_J$

In the finite-dimensional case where  $V = \mathbb{R}^k$ , it is well-known that the null space step  $\xi_J(x)$  defined by (1.4) is the orthogonal projection of the gradient  $\nabla J(x)$  onto the null space of the constraints

$$\text{Ker}(\mathbf{Dg}(x)) = \{\xi \in V \mid \mathbf{Dg}(x)\xi = 0\},$$

which is also the tangent space at  $x$  to the manifold  $\{y \in V \mid \mathbf{g}(y) = \mathbf{g}(x)\}$ . Of course, this is still true when  $V$  is a Hilbert space, as recalled in the next lemma.

**Lemma 1.** Let  $x \in V$  be a point satisfying the qualification condition (2.6). The following properties hold:

- (1) The space  $V$  has the following orthogonal decomposition:

$$V = \text{Ker}(\mathbf{Dg}(x)) \oplus \text{Ran}(\mathbf{Dg}(x)^T),$$

where we have introduced the range  $\text{Ran}(\mathbf{Dg}(x)^T) := \{\mathbf{Dg}(x)^T \lambda \mid \lambda \in \mathbb{R}^p\}$  of  $\mathbf{Dg}(x)^T$ .

Moreover, the operator  $\Pi_{g(x)} : V \rightarrow V$  defined by  $\Pi_{g(x)} = \mathbf{I} - \mathbf{Dg}^T(\mathbf{DgDg}^T)^{-1}\mathbf{Dg}(x)$  is the orthogonal projection onto  $\text{Ker}(\mathbf{Dg}(x))$ .

- (2) When  $\Pi_{g(x)}(\nabla J(x)) \neq 0$ ,  $-\xi_J(x) = -\Pi_{g(x)}(\nabla J(x))$  is the best normalized feasible descent direction for  $J$  in the sense that

$$-\frac{\xi_J(x)}{\|\xi_J(x)\|_V} = \arg \min_{\xi \in V} \mathbf{D}J(x)\xi \quad \text{s.t.} \begin{cases} \mathbf{Dg}(x)\xi = 0 \\ a(\xi, \xi) \leq 1. \end{cases} \quad (2.12)$$

- (3) The null space direction  $\xi_J(x) = \Pi_{g(x)}(\nabla J(x))$  is the closest least squares approximation to  $\nabla J(x)$  within the space  $\text{Ker}(\mathbf{Dg}(x))$ . It alternatively reads

$$\xi_J(x) = \nabla J(x) + \mathbf{Dg}(x)^T \lambda^*(x), \quad (2.13)$$

where the Lagrange multiplier  $\lambda^*(x) := -(\mathbf{DgDg}^T)^{-1}\mathbf{Dg}\nabla J(x)$  is the unique solution to the following least squares problem that is the dual of (2.12):

$$\lambda^*(x) = \arg \min_{\lambda \in \mathbb{R}^p} \|\nabla J(x) + \mathbf{Dg}(x)^T \lambda\|_V. \quad (2.14)$$

*Proof.*



- (1) Any  $\xi \in V$  may be decomposed as  $\xi = \Pi_{g(x)}(\xi) + (I - \Pi_{g(x)})(\xi)$ , where it is straightforward to verify that  $\Pi_{g(x)}(\xi) \in \text{Ker}(\text{Dg}(x))$ , and  $(I - \Pi_{g(x)})(\xi) \in \text{Ran}(\text{Dg}(x)^T)$ . In addition,  $\text{Ker}(\text{Dg}(x))$  and  $\text{Ran}(\text{Dg}(x)^T)$  are orthogonal for the inner product  $a$  since from (2.3), one has,

$$\forall \zeta \in \text{Ker}(\text{Dg}(x)), \forall \lambda \in \mathbb{R}^p, a(\text{Dg}(x)^T \lambda, \zeta) = \lambda^T \text{Dg}(x) \zeta = 0.$$

- (2) It follows from the first point that for any  $\xi \in \text{Ker}(\text{Dg}(x))$  such that  $\|\xi\|_V \leq 1$ ,

$$\text{DJ}(x)\xi = a(\nabla J(x), \xi) = a(\Pi_{g(x)}(\nabla J(x)), \xi) \geq -\|\Pi_{g(x)}(\nabla J(x))\|_V,$$

whence we easily infer that  $\xi := -\Pi_{g(x)}(\nabla J(x))/\|\Pi_{g(x)}(\nabla J(x))\|_V$  is the global minimizer of (2.12).

- (3) The Pythagore identity yields, for any  $\xi \in \text{Ker}(\text{Dg}(x))$ ,

$$\|\nabla J(x) - \xi\|_V^2 = \|(I - \Pi_{g(x)})\nabla J(x)\|_V^2 + \|\Pi_{g(x)}\nabla J(x) - \xi\|_V^2 \geq \|\nabla J(x) - \Pi_{g(x)}\nabla J(x)\|_V^2.$$

Hence the orthogonal projection  $\Pi_{g(x)}(\nabla J(x))$  is the best approximation of  $\nabla J(x)$  on  $\text{Ker}(\text{Dg}(x))$ . Recalling from the first point that  $\text{Ran}(\text{Dg}(x)^T)$  is the orthogonal complement of  $\text{Ker}(\text{Dg}(x))$ , we obtain also, for any  $\lambda \in \mathbb{R}^p$ ,

$$\|\Pi_{g(x)}(\nabla J(x))\|_V = \|\nabla J(x) - (I - \Pi_{g(x)})(\nabla J(x))\|_V \leq \|\nabla J(x) - \text{Dg}(x)^T \lambda\|_V,$$

whence the expression (2.13) and the minimization property (2.14) follow. Note that the uniqueness of the solution  $\lambda^*(x)$  to (2.14) results from the qualification condition (2.6).

Finally, the optimization problem (2.12) can be rewritten as

$$\min_{\substack{\xi \in V \\ a(\xi, \xi) \leq 1}} \max_{\lambda \in \mathbb{R}^p} \text{DJ}(x)\xi + \lambda^T \text{Dg}(x)\xi.$$

Hence the (formal) dual problem of (2.12) reads:

$$\max_{\lambda \in \mathbb{R}^p} \min_{\substack{\xi \in V \\ a(\xi, \xi) \leq 1}} \text{DJ}(x)\xi + \lambda^T \text{Dg}(x)\xi.$$

According to the definitions (2.3) and (2.4) of the gradient and of the Hilbertian transpose, the latter problem rewrites:

$$\max_{\lambda \in \mathbb{R}^p} \min_{\substack{\xi \in V \\ a(\xi, \xi) \leq 1}} a(\nabla J(x) + \text{Dg}(x)^T \lambda, \xi) = - \max_{\lambda \in \mathbb{R}^p} \|\nabla J + \text{Dg}^T \lambda\|_V,$$

where for given  $\lambda \in \mathbb{R}^p$ , the value

$$\xi^* := \frac{\nabla J(x) + \text{Dg}(x)^T \lambda}{\|\nabla J(x) + \text{Dg}(x)^T \lambda\|_V}$$

is that achieving the minimum in the minimization problem at the left-hand side of the above identity. This shows that (2.14) is the dual problem of (2.12). □

### 2.2.2. Properties of the range space step $\xi_C$

The next lemma is also classical in the literature. It characterizes the range space step  $\xi_C(x)$ , defined by (2.9), as the unique Gauss-Newton direction for the minimization of the constraint function  $g(x)$  which is orthogonal to the (linearized) set of constraints:

**Lemma 2.** *Let  $x \in V$  satisfy the condition (2.6); then:*

- (1) *The range space step  $\xi_C(x) = \text{Dg}^T(\text{DgDg}^T)^{-1}g(x)$  is orthogonal to  $\text{Ker}(\text{Dg}(x))$ :*

$$\forall \xi \in \text{Ker}(\text{Dg}(x)), a(\xi_C(x), \xi) = 0.$$

- (2)  *$-\xi_C(x)$  is a descent direction for the violation of the constraints:*

$$\text{Dg}(x)(-\xi_C(x)) = -g(x). \tag{2.15}$$

(3) The set of solutions to the Gauss-Newton program

$$\min_{\xi \in V} \|\mathbf{g}(x) + \mathbf{Dg}(x)\xi\|^2 \quad (2.16)$$

is the affine subspace  $\{-\xi_C(x) + \zeta \mid \zeta \in \text{Ker}(\mathbf{Dg}(x))\}$  of  $V$ .

*Proof.*

(1) This easily follows from (2.3) and the calculation:

$$\forall \xi \in \text{Ker}(\mathbf{Dg}(x)), a(\xi_C(x), \xi) = ((\mathbf{DgDg}^T)^{-1}\mathbf{g}(x))^T(\mathbf{Dg}(x)\xi) = 0.$$

(2) This is an immediate consequence of the definition (2.9) of  $\xi_C(x)$ . Note that (2.15) means that  $-\xi_C(x)$  is a descent direction for the violation of the constraints in the sense that it ensures that any coordinate  $g_i(x)$ ,  $i = 1, \dots, p$ , decreases along  $-\xi_C(x)$  if  $g_i(x) \geq 0$  and increases if  $g_i(x) \leq 0$ .

(3) Since (2.16) is a convex optimization problem, a necessary and sufficient condition for  $\xi \in V$  to be one solution is given by the usual first-order condition:

$$\forall \zeta \in V, (\mathbf{g}(x) + \mathbf{Dg}(x)\xi)^T(\mathbf{Dg}(x)\zeta) = a(\mathbf{Dg}(x)^T(\mathbf{g}(x) + \mathbf{Dg}(x)\xi), \zeta) = 0,$$

which rewrites:

$$\mathbf{Dg}(x)^T \mathbf{Dg}(x)\xi = -\mathbf{Dg}(x)^T \mathbf{g}(x).$$

Since the matrix  $(\mathbf{DgDg}^T)$  is invertible, this is in turn equivalent to:

$$\mathbf{Dg}(x)\xi = -\mathbf{g}(x).$$

Finally, (2) states that  $-\xi_C(x)$  is one particular solution to the above equation; therefore, any two solutions of this problem differ by some  $\zeta$  such that  $\mathbf{Dg}(x)\zeta = 0$ . □

### 2.3. Behavior of the trajectories of the flow

The main features of the definitions of  $\xi_J(x)$  and  $\xi_C(x)$  are the facts that  $\xi_J$  is orthogonal to the set of constraints, i.e.  $\mathbf{Dg}(x)\xi_J(x) = 0$ , and that  $-\xi_C(x)$  decreases the violation of the constraints while being orthogonal to  $\xi_J(x)$ . These ensure that the values of the constraint functional  $\mathbf{g}(x(t))$  decrease to zero along the trajectories of the ODE (1.3), independently of the value of  $\xi_J(x)$ . Then, as soon as the violation of the constraint becomes sufficiently small, the objective function  $J$  decreases without affecting the asymptotic vanishing of  $\mathbf{g}(x(t))$ . We review these properties in the next proposition, which was also observed in [48] in the finite-dimensional context.

**Proposition 1.** *Assume that the trajectories  $x(t)$  of the flow*

$$\begin{cases} \dot{x} = -\alpha_J(\mathbf{I} - \mathbf{Dg}^T(\mathbf{DgDg}^T)^{-1}\mathbf{Dg}(x))\nabla J(x) - \alpha_C\mathbf{Dg}^T(\mathbf{DgDg}^T)^{-1}\mathbf{g}(x) \\ x(0) = x_0 \end{cases} \quad (2.17)$$

*exist on some time interval  $[0, T]$  for  $T > 0$ , and that the qualification condition (2.6) holds at any point  $x(t)$ ,  $t \in [0, T]$ . Then the following properties hold true:*

(1) *The violation of the constraints decreases exponentially:*

$$\forall t \in [0, T], \mathbf{g}(x(t)) = e^{-\alpha_C t} \mathbf{g}(x_0). \quad (2.18)$$

(2) *Assume in addition that  $\text{rank}(\mathbf{Dg}) = p$  on  $K = \{x \in V \mid \|\mathbf{g}(x)\|_\infty \leq \|\mathbf{g}(x_0)\|_\infty\}$  and that*

$$\sup_{x \in K} \|\nabla J(x)\|_V |\sigma_p^{-1}(x)| < +\infty, \quad (2.19)$$

*where  $\sigma_p(x)$  is the smallest singular value of  $\mathbf{Dg}(x)$ . Then  $J$  decreases as long as the projection of its gradient on the set of constraints is large with respect to the decrease rate of the constraints given by (2.18), namely there exists a constant  $C > 0$  such that*

$$\forall t \in [0, T], \|\Pi_{g(x)}(\nabla J(x(t)))\|_V^2 > Ce^{-\alpha_C t} \Rightarrow \frac{d}{dt} J(x(t)) < 0. \quad (2.20)$$

(3) Any stationary point  $x^*$  of (2.17) satisfies the first-order KKT conditions (2.7) of the optimization program (2.1), that is:

$$\begin{cases} \mathbf{g}(x^*) = 0 \\ \exists \boldsymbol{\lambda}^* \in \mathbb{R}^p, \nabla J(x^*) + \mathbf{Dg}^\mathcal{T}(x^*)\boldsymbol{\lambda}^* = \Pi_{g(x^*)}(\nabla J(x^*)) = 0. \end{cases} \quad (2.21)$$

*Proof.*

(1) Using the definition (2.17), the decreasing property (2.15) together with the fact that  $\boldsymbol{\xi}_J(x)$  is orthogonal to  $\text{Ker}(\mathbf{Dg}(x))$ , we obtain:

$$\frac{d}{dt}(\mathbf{g}(x(t))) = -\alpha_C \mathbf{g}(x(t)),$$

whence (2.18) follows easily.

(2) Let us introduce the eigenvalue decomposition

$$\mathbf{Dg}(x)\mathbf{Dg}(x)^\mathcal{T} = \sum_{i=1}^p \sigma_i(x)^2 \mathbf{u}_i(x)\mathbf{u}_i(x)^\mathcal{T}, \text{ where } \sigma_1(x) \geq \dots \geq \sigma_p(x) > 0, \mathbf{u}_i(x)^\mathcal{T}\mathbf{u}_j(x) = \delta_{ij},$$

of the symmetric, positive definite  $p \times p$  matrix  $\mathbf{Dg}(x)\mathbf{Dg}(x)^\mathcal{T}$ . Let then  $\mathbf{v}_i(x)^\dagger : V \rightarrow \mathbb{R}$  be the linear form defined for any  $\boldsymbol{\xi} \in V$  by  $\mathbf{v}_i(x)^\dagger \boldsymbol{\xi} = \sigma_i(x)^{-1} \mathbf{u}_i(x)^\mathcal{T} \mathbf{Dg}(x) \boldsymbol{\xi}$  and  $\mathbf{v}_i(x)$  be the vector in  $V$  such that  $\forall \boldsymbol{\xi} \in V, a(\mathbf{v}_i(x), \boldsymbol{\xi}) = \mathbf{v}_i(x)^\dagger \boldsymbol{\xi}$ ; more explicitly,  $\mathbf{v}_i(x) = \sigma_i(x)^{-1} \mathbf{Dg}(x)^\mathcal{T} \mathbf{u}_i(x)$ . These definitions allow to write a singular value decomposition for  $\mathbf{Dg}(x)$ ; it is indeed easily verified from the definitions of  $\mathbf{u}_i(x)$  and  $\mathbf{v}_i(x)$  that:

$$\mathbf{Dg}(x) = \sum_{i=1}^p \sigma_i(x) \mathbf{u}_i(x) \mathbf{v}_i(x)^\dagger, \text{ and } \mathbf{Dg}(x)^\mathcal{T} = \sum_{i=1}^p \sigma_i(x) \mathbf{v}_i(x) \mathbf{u}_i(x)^\mathcal{T}$$

with  $a(\mathbf{v}_i(x), \mathbf{v}_j(x)) = \mathbf{v}_i(x)^\dagger \mathbf{v}_j(x) = \delta_{ij}$ . We now calculate:

$$\mathbf{Dg}^\mathcal{T}(\mathbf{Dg}\mathbf{Dg}^\mathcal{T})^{-1} \mathbf{g}(x) = \sum_{i=1}^p \sigma_i^{-1}(x) (\mathbf{u}_i(x)^\mathcal{T} \mathbf{g}(x)) \mathbf{v}_i(x),$$

whence we obtain the following inequality:

$$\forall x \in V, |\mathbf{Dg}(x)\mathbf{Dg}^\mathcal{T}(\mathbf{Dg}\mathbf{Dg}^\mathcal{T})^{-1} \mathbf{g}| \leq \sigma_p^{-1}(x) \|\nabla J(x)\|_V \|\mathbf{g}(x)\|. \quad (2.22)$$

Since

$$\frac{d}{dt} J(x(t)) = -\alpha_J \mathbf{Dg}(x(t)) \boldsymbol{\xi}_J(x(t)) - \alpha_C \mathbf{Dg}(x(t)) \boldsymbol{\xi}_C(x(t)),$$

it follows that  $\frac{d}{dt} J(x(t)) < 0$  as soon as  $\alpha_J |\mathbf{Dg}(x(t)) \boldsymbol{\xi}_J(x(t))| > \alpha_C |\mathbf{Dg}(x(t)) \boldsymbol{\xi}_C(x(t))|$ . Thus, from (2.18) and (2.22), the constant  $C$  in (2.20) can be selected as

$$C = p \frac{\alpha_C}{\alpha_J} \|\mathbf{g}(x_0)\| \sup_{x \in K} [\sigma_p^{-1}(x) \|\nabla J(x)\|]. \quad (2.23)$$

(3) Since the vectors  $\boldsymbol{\xi}_J(x)$  and  $\boldsymbol{\xi}_C(x)$  are orthogonal for any point  $x \in V$ , a stationary point  $x^*$  of (2.17) must satisfy

$$\Pi_{g(x^*)}(\nabla J(x^*)) = 0, \text{ and } \mathbf{Dg}^\mathcal{T}(\mathbf{Dg}\mathbf{Dg}^\mathcal{T})^{-1} \mathbf{g}(x^*) = 0, \quad (2.24)$$

and so the first KKT condition in (2.7) is satisfied with the value  $\boldsymbol{\lambda} = -(\mathbf{Dg}\mathbf{Dg}^\mathcal{T})^{-1} \mathbf{Dg}(x^*) \nabla J(x^*)$  of the Lagrange multiplier. Then left multiplication by  $\mathbf{Dg}$  in the second identity in (2.24) implies  $\mathbf{g}(x^*) = 0$ , which completes the proof.  $\square$

*Remark 2.* The solutions to the dynamical system (2.17) are defined for small times if  $\boldsymbol{\xi}_J$  and  $\boldsymbol{\xi}_C$  are locally Lipschitz vector fields, which is the case if e.g.  $J$  and  $\mathbf{g}$  are of class  $\mathcal{C}^2$  [22]. In the case where  $V$  is finite-dimensional, the assumption (2.19) is satisfied if the set  $K = \{x \in V | \mathbf{g}(x) \leq \mathbf{g}(x_0)\}$  is bounded and the functions  $J$  and  $\mathbf{g}$  are  $\mathcal{C}^1$  functions. It is worth noting that even if not enough regularity assumptions hold to

ensure the existence of solutions to (2.17), similar properties to those of Proposition 1 hold for the discretized scheme

$$x_{n+1} = x_n - \Delta t(\alpha_J \xi_J(x_n) + \alpha_C \xi_C(x_n)), \quad (2.25)$$

which is sufficient for optimization. One can indeed verify that:

- (1) At first order, the constraints decrease with a geometric rate:  $\mathbf{g}(x_{n+1}) = (1 - \alpha_C \Delta t) \mathbf{g}(x_n) + o(\Delta t)$ .
- (2)  $x^*$  is an accumulation point of the sequence  $(x_n)_{n \in \mathbb{N}}$  if and only if  $\mathbf{g}(x^*) = 0$  and  $x^*$  is a KKT point of the problem (2.1), satisfying (2.7).

*Remark 3.* In our design of the update rule (1.3) to (1.5), it is possible to control more accurately the pace at which each of the constraints decreases: consider a diagonal matrix of positive coefficients  $\mathbf{K} = \text{diag}(\kappa_i)_{1 \leq i \leq p}$  and replace the definition (1.5) or (2.9) of  $\xi_C(x)$  by

$$\xi_C(x) := \mathbf{D} \mathbf{g}^T (\mathbf{D} \mathbf{g} \mathbf{D} \mathbf{g}^T)^{-1} \mathbf{K} \mathbf{g}(x).$$

Then it can be shown along the lines of the previous discussion that each constraint function  $g_i$  decreases at its own rate  $\kappa_i \alpha_C$  along the solution  $x(t)$  of (2.17):

$$\forall t \in [0, T], g_i(x(t)) = e^{-\kappa_i \alpha_C t} g_i(x_0).$$

### 3. PROPOSED EXTENSION TO EQUALITY AND INEQUALITY CONSTRAINTS

We now proceed to extend the dynamical system (1.3) to handle inequality constraints as well. We consider from now on the full optimization problem

$$\begin{aligned} \min_{x \in \mathcal{X}} \quad & J(x) \\ \text{s.t.} \quad & \begin{cases} \mathbf{g}(x) = 0 \\ \mathbf{h}(x) \leq 0, \end{cases} \end{aligned} \quad (3.1)$$

taking place over the Hilbert space  $\mathcal{X} = V$  with inner product  $a(\cdot, \cdot)$ , and where  $J : V \rightarrow \mathbb{R}$ ,  $\mathbf{g} : V \rightarrow \mathbb{R}^p$  and  $\mathbf{h} : V \rightarrow \mathbb{R}^q$  are differentiable functions.

This section is organized as follows: after setting notations in Section 3.1, we briefly review in Section 3.2 the classical method of slack variables for handling the inequality constraints in the problem (3.1). Section 3.3 then introduces another method for dealing with these constraints, which is original to the best of our knowledge. The essence of our method is the resolution of a dual quadratic subproblem for identifying the subset of constraints whose violation ‘naturally’ decreases in the course of the minimization of  $J$ , and those which should be enforced by projection of the gradient of  $J$ . The behavior of the induced flow in the context of the problem (3.1) is finally analyzed in Section 3.3.3 in the spirit of Proposition 1.

#### 3.1. Notations and preliminaries

The set of indices corresponding to saturated or violated inequality constraints at  $x \in V$  is denoted by  $\tilde{I}(x)$

$$\tilde{I}(x) = \{i \in \{1, \dots, q\} \mid h_i(x) \geq 0\}, \quad (3.2)$$

and  $\tilde{q}(x) := \text{Card}(\tilde{I}(x))$  is the number of such constraints. Recall the notation  $\mathbf{h}_I(x) = (h_i(x))_{i \in I}$  for the inequality constraints indexed by a subset  $I \subset \tilde{I}(x)$  and  $\mathbf{C}_I(x)$ , defined by (1.6), for the vector collecting the equality constraints  $\mathbf{g}(x)$  and those selected inequality constraints  $\mathbf{h}_I(x)$ . The constraints are said to be qualified at  $x \in V$ , in the sense that the linearized saturated or violated constraints are independent:

$$\text{rank}(\mathbf{D} \mathbf{C}_{\tilde{I}(x)}(x)) = p + \tilde{q}(x). \quad (3.3)$$

If the point  $x$  satisfies the constraints, (3.3) is one usual qualification condition (of course, there are other possible qualification conditions, see [14, 36]). Define  $\Pi_{\mathbf{C}_I} : V \rightarrow V$ , the orthogonal projection operator onto  $\text{Ker}(\mathbf{D} \mathbf{C}_I(x))$ , by

$$\Pi_{\mathbf{C}_I} = I - \mathbf{D} \mathbf{C}_I(x)^T (\mathbf{D} \mathbf{C}_I(x) \mathbf{D} \mathbf{C}_I(x)^T)^{-1} \mathbf{D} \mathbf{C}_I(x), \quad (3.4)$$

and  $(\boldsymbol{\lambda}_I(x), \boldsymbol{\mu}_I(x)) \in \mathbb{R}^p \times \mathbb{R}_+^{\text{Card}(I)}$  the corresponding Lagrange multipliers:

$$\begin{bmatrix} \boldsymbol{\lambda}_I(x) \\ \boldsymbol{\mu}_I(x) \end{bmatrix} := -(\text{DC}_I \text{DC}_I^T)^{-1} \text{DC}_I(x) \nabla J(x). \quad (3.5)$$

Last but not least, let us recall that in the present context of equality and inequality constrained problem (3.1), the necessary first-order optimality conditions (the KKT conditions) for a given point  $x^* \in V$ , satisfying the qualification condition (3.3), read as follows: there exist  $\boldsymbol{\lambda}(x^*) \in \mathbb{R}^p$  and  $\boldsymbol{\mu}(x^*) \in \mathbb{R}_+^q$  such that

$$\begin{cases} \nabla J(x^*) + \text{D}\mathbf{g}(x^*)^T \boldsymbol{\lambda}(x^*) + \text{D}\mathbf{h}(x^*)^T \boldsymbol{\mu}(x^*) = 0, \\ \mathbf{g}(x^*) = 0, \quad \mathbf{h}(x^*) \leq 0, \\ \forall i = 1, \dots, q, \quad \mu_i h_i(x^*) = 0; \end{cases} \quad (3.6)$$

see again [14, 36].

### 3.2. The method of slack variables for inequality constraints

It is classical to introduce slack variables so as to turn inequality constraints into equality constraints of an augmented problem including these additional variables, see [38] for the present gradient flow in the finite dimensional context. In other words, problem (3.1) is replaced with the following equivalent one, involving  $q$  extra variables  $(z_1, \dots, z_q) \in \mathbb{R}^q$ :

$$\begin{aligned} \min_{\substack{x \in V \\ z \in \mathbb{R}^q}} J(x) \\ \text{s.t.} \quad \mathbf{C}(x, z) = 0, \end{aligned} \quad (3.7)$$

where the augmented vector of constraints  $\mathbf{C}(x, z)$  reads:

$$\mathbf{C}(x, z) := \begin{bmatrix} \mathbf{g}(x) \\ h_1(x) + \frac{1}{2} z_1^2 \\ \vdots \\ h_q(x) + \frac{1}{2} z_q^2 \end{bmatrix} \in \mathbb{R}^{p+q}.$$

Problem (3.7) is an equality constrained optimization problem of the form (2.1), set over the Hilbert space

$$\tilde{V} := V \times \mathbb{R}^q, \text{ with inner product } \tilde{a}((x, z), (x', z')) := a(x, x') + z^T z'.$$

It can be solved thanks to the proposed algorithm in Section 2; the associated gradient flow for (3.7) reads:

$$\begin{cases} \begin{bmatrix} \dot{x}(t) \\ \dot{z}(t) \end{bmatrix} = -\alpha_J (\text{I} - \text{DC}^T (\text{DCDC}^T)^{-1} \text{DC}) \begin{bmatrix} \nabla J(x(t)) \\ 0 \end{bmatrix} - \alpha_C \text{DC}^T (\text{DCDC}^T)^{-1} \mathbf{C}(x(t), z(t)), \\ \begin{bmatrix} x(0) \\ z(0) \end{bmatrix} = \begin{bmatrix} x_0 \\ z_0 \end{bmatrix}, \end{cases} \quad (3.8)$$

where  $x_0 \in V$  is the considered initial point in the resolution of (3.1), and the variable  $z$  is initialized with a value  $z_0 \in \mathbb{R}^q$  in such a way that the inequality constraints of (3.1) which are inactive for  $x_0$  (i.e.  $h_i(x_0) < 0$ ) are associated with satisfied equality constraints  $C_{p+i}(x_0, z_0) = 0$  in (3.7):

$$\forall i \in \{1, \dots, q\}, \quad z_{0,i} = \sqrt{2|h_i(x_0)|}.$$

In the finite-dimensional setting  $V = \mathbb{R}^k$  and when  $J$ ,  $\mathbf{g}$  and  $\mathbf{h}$  are  $\mathcal{C}^2$  functions, Schropp and Singer proved in [38] that:

- (i) Stationary points of the extended flow (3.8) are exactly critical points of (3.1), that is points  $x^*$  satisfying (3.6) but with  $\boldsymbol{\mu}(x^*) \in \mathbb{R}^q$  possibly negative.

- (ii) Among all possible critical points, only KKT points (fulfilling all three conditions (3.6) with  $\mu(x^*) \in \mathbb{R}_+^q$ ) are asymptotically stable equilibria.

As a consequence, the solution vector  $x(t)$  to (3.8) converges in practice to a KKT point for problem (3.1); see also Section 3.4 about this point.

In the following, we shall introduce a different approach that does not use slack variables, but rather relies on a suitable selection of the active or violated constraints. The main differences with the method of slack variables lies in that inactive constraints are ignored, and that all stationary points of our proposed gradient flow are KKT points (satisfying all conditions (3.6)).

### 3.3. The proposed algorithm

Inspired by the methodology developed in Section 2, we still propose to solve the equality and inequality constrained problem (3.1) thanks to a dynamical system of the form:

$$\begin{cases} \dot{x}(t) = -\alpha_J \xi_J(x(t)) - \alpha_C \xi_C(x(t)) \\ x(0) = x_0, \end{cases} \quad (3.9)$$

whose discretized version reads:

$$x_{n+1} = x_n - \Delta t (\alpha_J \xi_J(x_n) + \alpha_C \xi_C(x_n)). \quad (3.10)$$

In the next subsections, we define  $\xi_J(x)$  and  $\xi_C(x)$  from formulas analogous to (1.7) and (1.8), which involve a procedure discriminating a relevant subset  $\hat{I}(x) \subset \tilde{I}(x)$  of the saturated or violated constraints. Finally, we establish the properties of the flow (3.9) in Proposition 5.

#### 3.3.1. Definition of the range step direction

**Definition 3.** For the optimization problem (3.1), the range step  $\xi_C(x)$  is defined by

$$\xi_C(x) := DC_{\tilde{I}(x)}^T (DC_{\tilde{I}(x)} DC_{\tilde{I}(x)}^T)^{-1} C_{\tilde{I}(x)}(x), \quad (3.11)$$

where  $\tilde{I}(x)$  is the subset of saturated or violated constraints, defined by (3.2).

The purpose of the range space step  $\xi_C(x)$  is to decrease the violation of the constraints as we shall see in Proposition 5 below. The counterpart of Lemma 2 holds exactly in this context, in particular:

- (1)  $\xi_C(x)$  is orthogonal to  $\text{Ker}(DC_{\tilde{I}(x)})$ .
- (2)  $-\xi_C(x)$  is a Gauss-Newton direction for the violation of the constraints:

$$DC_{\tilde{I}(x)}(-\xi_C(x)) = -C_{\tilde{I}(x)}(x).$$

The definition of the null space direction  $\xi_J(x)$  is slightly more involved as it is not obtained by replacing  $Dg(x)$  by  $DC_{\tilde{I}(x)}$  in (2.8). It requires the introduction of a different subset  $\hat{I}(x) \subset \tilde{I}(x)$ , which is now detailed.

#### 3.3.2. Definition and characterizations of the null space direction $\xi_J(x)$

Inspired by the characterization of the null space direction of Lemma 1 for equality constrained problems, the null space direction  $\xi_J(x)$  is now sought, up to a change of sign, as a best normalized descent direction diminishing violated or saturated inequality constraints, i.e.  $-\xi_J(x)$  shall be set positively proportional to the solution of the following minimization problem:

$$\begin{aligned} \min_{\xi \in V} \quad & DJ(x)\xi \\ \text{s.t.} \quad & \begin{cases} Dg(x)\xi = 0 \\ D\mathbf{h}_{\tilde{I}(x)}(x)\xi \leq 0 \\ \|\xi\|_V \leq 1. \end{cases} \end{aligned} \quad (3.12)$$

The purpose of this subsection is to characterize explicitly the minimizer  $\xi^*(x)$  of (3.12) by relying on its dual problem. As a consequence an explicit formula for the null space direction  $\xi_J(x)$ , in the form of (1.7), will be given in Definition 4 below.

We now introduce the dual optimization problem to (3.12), which is analogous to the dual problem (2.14) of the previous section.

**Proposition 2.** *Let  $x \in V$  satisfy the qualification condition (3.3). There exists a unique couple of multipliers  $\lambda^*(x) \in \mathbb{R}^p$  and  $\mu^*(x) \in \mathbb{R}_+^{\tilde{q}(x)}$  solution to the following quadratic optimization problem which is the dual of (3.12):*

$$(\lambda^*(x), \mu^*(x)) := \arg \min_{\substack{\lambda \in \mathbb{R}^p \\ \mu \in \mathbb{R}_+^{\tilde{q}(x)}, \mu \geq 0}} \|\nabla J(x) + Dg(x)^\top \lambda + Dh_{\tilde{I}(x)}(x)^\top \mu\|_V. \quad (3.13)$$

*Proof.* Problem (3.12) is equivalent to the following min-max formulation:

$$\min_{\substack{\xi \in V \\ a(\xi, \xi) \leq 1}} \max_{\substack{\lambda \in \mathbb{R}^p \\ \mu \in \mathbb{R}_+^{\tilde{q}(x)}}} DJ(x)\xi + \lambda^\top Dg(x)\xi + \mu^\top Dh_{\tilde{I}(x)}(x)\xi.$$

Inverting formally the min and the max and performing the maximization with respect to  $\xi$  as in the proof of Lemma 1 yields that (3.13) is the dual problem of (3.12) up to a change of sign (the duality gap between (3.13) and (3.12) will be shown to vanish in Proposition 3). The program (3.13) brings into play the closed convex set  $\mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x)}$  and the least squares functional

$$(\lambda, \mu) \mapsto \left\| \nabla J(x) + DC_{\tilde{I}(x)}(x)^\top \begin{bmatrix} \lambda \\ \mu \end{bmatrix} \right\|_V.$$

The latter is strictly convex over  $\mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x)}$  by virtue of (3.3). Hence, (3.13) has a unique solution.  $\square$

The optimization problem (3.13) belongs to the class of non negative least squares problems; it can be solved efficiently with a number of dedicated solvers, such as `cvxopt` [8] or `IPOPT` [46]. One nice feature of (3.13) lies in that its dimension is the number  $p + \tilde{q}(x)$  of saturated or violated constraints, which can be small for many practical cases, as e.g. in our shape optimization applications of Section 6. It is also possible to exploit the sparsity of the constraints if  $p + \tilde{q}(x)$  is large, see Remark 5 below.

In the next proposition, we relate the optimal values and the solutions  $\xi^*(x)$  and  $(\lambda^*(x), \mu^*(x))$  of the primal and dual problems (3.12) and (3.13). In essence, we show that the optimal feasible descent direction  $\xi^*(x)$  of (3.15) is the projection of the gradient  $\nabla J(x)$  onto the cone of feasible directions. The proof follows classical arguments of linear programming duality theory and it is detailed for the convenience of the reader.

**Proposition 3.** *Let  $x \in V$  satisfy the qualification condition (3.3) and denote*

$$m^*(x) := \|\nabla J(x) + Dg(x)^\top \lambda^*(x) + Dh_{\tilde{I}(x)}(x)^\top \mu^*(x)\|_V$$

*the value of the dual problem (3.13). Then the value of the primal problem (3.12) is  $p^*(x) = -m^*(x)$  and the following alternative holds:*

- (1)  $m^*(x) = 0$ : the first line of the KKT conditions (3.6) for the minimization problem (3.1) holds with (necessarily unique) Lagrange multipliers  $(\lambda^*(x), \mu^*(x)) \in \mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x)}$ :

$$\nabla J(x) + Dg(x)^\top \lambda^*(x) + Dh_{\tilde{I}(x)}(x)^\top \mu^*(x) = 0. \quad (3.14)$$

*One particular minimizer of (3.12) is  $\xi^*(x) = 0$ .*

- (2)  $m^*(x) > 0$ : (3.14) does not hold and there exists a unique minimizer  $\xi^*(x)$  to (3.12), given by

$$\xi^*(x) = - \frac{\nabla J(x) + Dg(x)^\top \lambda^*(x) + Dh_{\tilde{I}(x)}(x)^\top \mu^*(x)}{\|\nabla J(x) + Dg(x)^\top \lambda^*(x) + Dh_{\tilde{I}(x)}(x)^\top \mu^*(x)\|_V}. \quad (3.15)$$



*Proof.* Let  $\xi \in V$  be a feasible direction for the problem (3.12), i.e.  $Dg(x)\xi = 0$ ,  $Dh_{\tilde{I}(x)}(x)\xi \leq 0$  and  $\|\xi\|_V \leq 1$ . Then for any  $(\lambda, \mu) \in \mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x)}$ , it holds

$$\begin{aligned} DJ(x)\xi &\geq DJ(x)\xi + \lambda^T Dg(x)\xi + \mu^T Dh_{\tilde{I}(x)}(x)\xi \\ &= a(\nabla J(x) + Dg(x)^T \lambda + Dh_{\tilde{I}(x)}(x)^T \mu, \xi) \\ &\geq -\|\nabla J(x) + Dg(x)^T \lambda + Dh_{\tilde{I}(x)}(x)^T \mu\|_V \end{aligned} \quad (3.16)$$

Since (3.16) holds for any feasible direction  $\xi$  for (3.12), and for any  $(\lambda, \mu) \in \mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x)}$ , it follows:

$$\min_{\substack{\xi \in V \\ \xi \text{ feasible for (3.12)}}} DJ(x)\xi \geq - \min_{\substack{\lambda \in \mathbb{R}^p \\ \mu \in \mathbb{R}_+^{\tilde{q}(x)}, \mu \geq 0}} \|\nabla J(x) + Dg(x)^T \lambda + Dh_{\tilde{I}(x)}(x)^T \mu\|_V. \quad (3.17)$$

Therefore, we have proven that  $p^*(x) \geq -m^*(x)$ . We now examine the alternative  $m^*(x) = 0$  or  $m^*(x) > 0$ :

- (1) If  $m^*(x) = 0$ , then (3.17) implies  $p^*(x) \geq 0$ . Therefore, the value of (3.12) is  $p^*(x) = -m^*(x) = 0$ , attained in particular at  $\xi^* = 0$ , and more generally at any feasible  $\xi^* \in V$  satisfying  $\mu^*(x)^T Dh_{\tilde{I}(x)}(x)\xi^* = 0$ , as follows readily from the KKT conditions for (3.12). Furthermore, the KKT equation (3.14) is satisfied by definition of  $m^*(x) = 0$ .
- (2) Assume now  $m^*(x) > 0$ . The KKT condition for (3.12) states that for any local optimum  $\xi'$ , there exists  $(\lambda', \mu') \in \mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x)}$  and  $\alpha \geq 0$  such that,

$$\forall \xi \in V, (DJ(x) + \lambda'^T Dg(x) + \mu'^T Dh_{\tilde{I}(x)}(x))\xi = -\alpha a(\xi', \xi). \quad (3.18)$$

Using Riesz identifications of the gradient and the differentials, we obtain

$$\alpha \xi' = -(\nabla J(x) + Dg(x)^T \lambda' + Dh_{\tilde{I}(x)}(x)^T \mu'),$$

and since  $m^*(x) > 0$ , it is necessary that  $\alpha > 0$ . The complementarity condition  $\alpha(a(\xi', \xi') - 1) = 0$  yields then  $\|\xi'\|_V = 1$ , which readily implies:

$$\xi' = -\frac{\nabla J(x) + Dg(x)^T \lambda' + Dh_{\tilde{I}(x)}(x)^T \mu'}{\|\nabla J(x) + Dg(x)^T \lambda' + Dh_{\tilde{I}(x)}(x)^T \mu'\|_V}.$$

Then the complementarity condition for (3.12) implies  $\mu'^T Dh_{\tilde{I}(x)}(x)\xi' = 0$ . Therefore it holds that

$$\begin{aligned} DJ(x)\xi' &= DJ(x)\xi' + \lambda'^T Dg(x)\xi' + \mu'^T Dh_{\tilde{I}(x)}(x)\xi' \\ &= a(\nabla J(x) + Dg(x)^T \lambda' + Dh_{\tilde{I}(x)}(x)^T \mu', \xi') \\ &= -\|\nabla J(x) + Dg(x)^T \lambda' + Dh_{\tilde{I}(x)}(x)^T \mu'\|_V. \end{aligned} \quad (3.19)$$

The previous equation together with the inequality (3.16) with  $\xi = \xi'$  then implies that  $(\lambda', \mu')$  achieves the minimum of (3.13). By uniqueness, this implies  $\lambda' = \lambda^*(x)$  and  $\mu' = \mu^*(x)$ , hence  $\xi' = \xi^*(x)$ . Furthermore,  $p^*(x) = DJ(x)\xi^*(x) = DJ(x)\xi' = -m^*(x)$ .

□

Finally, the next proposition characterizes explicitly the expression of the optimal descent direction  $\xi^*(x)$  from the signs of the multiplier  $\mu^*(x)$ , and highlights in which sense the problem (3.12) is combinatorial. Recall definitions (3.4) and (3.5) for the projection operator  $\Pi_{C_I}$  and the multipliers  $(\lambda_I(x), \mu_I(x))$ .

**Proposition 4.** *In the context of point (2) in Proposition 3, let  $\xi^*(x)$  and  $(\lambda^*(x), \mu^*(x))$  be the minimizers of the primal and dual problems (3.12) and (3.13). Define the subset  $\hat{I}(x) \subset \tilde{I}(x)$  by*

$$\hat{I}(x) := \{i \in \tilde{I}(x) \mid \mu_i^*(x) > 0\}. \quad (3.20)$$

(1)  $(\lambda^*(x), \mu^*(x))$  and  $\xi^*(x)$  are explicitly given in terms of  $\hat{I}(x)$  by:

$$\begin{bmatrix} \lambda^*(x) \\ \hat{\mu}^*(x) \end{bmatrix} = \begin{bmatrix} \lambda_{\hat{I}(x)}(x) \\ \mu_{\hat{I}(x)}(x) \end{bmatrix} = -(\text{DC}_{\hat{I}(x)} \text{DC}_{\hat{I}(x)}^T)^{-1} \text{DC}_{\hat{I}(x)} \nabla J(x), \quad (3.21)$$

$$\xi^*(x) = -\frac{\Pi_{C_{\hat{I}(x)}}(\nabla J(x))}{\|\Pi_{C_{\hat{I}(x)}}(\nabla J(x))\|_V}, \quad (3.22)$$

where  $\hat{\mu}^*(x) := (\mu_i^*(x))_{i \in \hat{I}(x)}$  is the vector collecting all positive components of  $\mu^*(x)$ .

(2)  $\hat{I}(x)$  is equivalently the unique solution to either of the following discrete optimization problems:

$$\begin{aligned} \hat{I}(x) &= \arg \max_{I \subset \tilde{I}(x)} \|\Pi_{C_I}(\nabla J(x))\|_V \\ \text{s.t. } & \text{Dh}_{\tilde{I}(x)}(x) \Pi_{C_I}(\nabla J(x)) \geq 0, \end{aligned} \quad (3.23)$$

$$\begin{aligned} \hat{I}(x) &= \arg \min_{I \subset \tilde{I}(x)} \|\Pi_{C_I}(\nabla J(x))\|_V \\ \text{s.t. } & \mu_I(x) \geq 0. \end{aligned} \quad (3.24)$$

In particular,  $\hat{I}(x)$  is the unique subset  $I \subset \tilde{I}(x)$  satisfying simultaneously both feasibility conditions

$$\text{Dh}_{\tilde{I}(x)}(x) \Pi_{C_I}(\nabla J(x)) \geq 0 \text{ and } \mu_I(x) \geq 0.$$

*Proof.*

(1) The complementary condition for the primal and dual problems (3.12) and (3.13) reads

$$\forall i \in \tilde{I}(x), \quad \mu_i^*(x) \text{Dh}_i(x) \xi^*(x) = 0. \quad (3.25)$$

Therefore,  $\text{Dh}_i(x) \xi^*(x) = 0$  for all indices  $i \in \hat{I}(x)$ , which implies that as  $\text{DC}_{\hat{I}(x)}(x) \xi^*(x) = 0$ . Then, after left multiplication of (3.15) by  $(\text{DC}_{\hat{I}(x)} \text{DC}_{\hat{I}(x)}^T)^{-1} \text{DC}_{\hat{I}(x)}$ , we obtain (3.21), whence (3.22) follows.

(2) Let  $I \subset \tilde{I}(x)$  a subset satisfying  $\text{Dh}_{\tilde{I}(x)}(x) \Pi_{C_I}(\nabla J(x)) \geq 0$ . This implies that

$$\xi = -\Pi_{C_I}(\nabla J(x)) / \|\Pi_{C_I}(\nabla J(x))\|_V$$

is feasible for the primal problem (3.12), and we obtain by definition of  $\xi^*(x)$  that

$$-\|\Pi_{C_{\hat{I}(x)}}(\nabla J(x))\|_V = \text{DJ}(x) \xi^*(x) \leq \text{DJ}(x) \xi = -\|\Pi_{C_I}(\nabla J(x))\|_V, \quad (3.26)$$

whence the maximization property (3.23).

For  $I \subset \tilde{I}(x)$  satisfying  $\mu_I(x) \geq 0$ , we obtain feasible multipliers  $(\lambda, \mu)$  for the dual problem (3.13) by taking  $\mu$  to be equal to  $\mu_I$  on the indices of  $I$  and extended by 0 in the complementary subset  $\tilde{I}(x) \setminus I$ . Then the optimality of  $(\lambda^*(x), \mu^*(x))$  for this dual problem reads:

$$\begin{aligned} \|\Pi_{C_{\hat{I}(x)}}(\nabla J(x))\|_V &= \|\nabla J + \text{Dg}(x)^T \lambda^*(x) + \text{Dh}_{\tilde{I}(x)}(x)^T \mu^*(x)\|_V \\ &\leq \|\nabla J(x) + \text{Dg}(x)^T \lambda + \text{Dh}_{\tilde{I}(x)}^T \mu\|_V = \|\Pi_{C_I}(\nabla J(x))\|_V, \end{aligned} \quad (3.27)$$

whence the minimization property (3.24).

□

In view of (3.20), the optimal multiplier  $\mu^*(x)$  can be interpreted as an indicator variable specifying which constraints of  $\tilde{I}(x)$  are ‘not aligned’ with the gradient  $\nabla J(x)$  and should be kept in the subset  $\hat{I}(x)$ . The best descent direction (in the sense of (3.12)) is obtained by projecting the gradient  $\nabla J(x)$  onto the tangent space of the constraint subset  $\hat{I}(x)$  rather than onto the full set of violated or saturated constraints  $\tilde{I}(x)$ . Indeed, the descent direction  $\xi = -\Pi_{C_{\tilde{I}(x)}} \nabla J(x)$  that would be obtained by projecting  $\nabla J(x)$  on the whole set  $\tilde{I}(x)$  would only keep them all constant at first order, i.e.  $\text{Dh}_i(x) \xi = 0$ , (see Remark 6 for more details). It is therefore more efficient to project  $\nabla J(x)$  only on those constraints associated to the indices  $i \in \hat{I}(x)$ ,

thus allowing the remaining ones (associated to  $i \in \tilde{I}(x) \setminus \hat{I}(x)$ , indicating vanishing multipliers  $\mu_i^*(x) = 0$ ) to decrease since the calculated descent direction ensures that  $Dh_i(x)\xi^*(x) \leq 0$  holds for all  $i = 1, \dots, q$ .

Note that actually, the use of a dual problem such as (3.13) in order to obtain information about which constraints should remain active is classical in active sets methods, see e.g. [16, 31, 36].

In principle, the subset  $\hat{I}(x)$  could be found by solving the discrete problems (3.23) or (3.24). However, we expect that in practice, it is more efficient to rely on iterative solvers relying on gradient descents for solving the dual problem (3.13), e.g. a cone programming solver or a non negative least squares algorithm such as [16]. This is what we do in the sequel.

*Remark 4.* With our notations, the optimization scheme proposed by Barbarosie et. al. [11, 12] reads

$$\begin{cases} x_{n+1} = x_n - \Delta t \nabla J(x_n) - DC_{I(x_n)}^\top \nu_n \\ \nu_n = -\Delta t (DC_{I(x_n)} DC_{I(x_n)}^\top)^{-1} DC_{I(x_n)} \nabla J(x_n) + DC_{I(x_n)}^\top (DC_{I(x_n)} DC_{I(x_n)}^\top)^{-1} C_{I(x_n)}, \end{cases} \quad (3.28)$$

where the set  $I(x_n)$  is obtained by removing indices from  $\tilde{I}(x_n)$  one by one, starting from the index  $i_0$  associated with the most negative multiplier  $\nu_{n,i_0} < 0$ , until all of them becomes non negative. Therefore, the set  $I(x_n)$  used in this strategy and that  $\hat{I}(x_n)$  featured in our strategy, given by (3.20), do not coincide in general; one could think of configurations where the procedure of [12] would fail to find the optimal set  $\hat{I}(x_n)$  (for example if  $i_0 \in \hat{I}(x_n)$ ) and would project the gradient on a less optimal subset of constraints. We note that no convergence result is given by the authors about this procedure.

Having introduced the subset  $\hat{I}(x)$  (defined in (3.20)), we are now able to define the null space direction  $\xi_J(x)$  in the present context:  $-\xi_J(x)$  is set to be a positive multiple of the optimal descent direction  $\xi^*(x)$  supplied by (3.22).

**Definition 4.** For any point  $x \in V$  satisfying the constraint qualification (3.3), the null space direction  $\xi_J(x)$  at  $x$  for the optimization problem (3.1) is defined by:

$$\xi_J(x) := \Pi_{C_{\hat{I}(x)}}(\nabla J(x)) = (I - DC_{\hat{I}(x)}(x)^\top (DC_{\hat{I}(x)} DC_{\hat{I}(x)}^\top)^{-1} DC_{\hat{I}(x)}) \nabla J(x), \quad (3.29)$$

where  $\hat{I}(x)$  is the set defined by (3.20).

The main point in Definition 4 is that, while all violated and saturated constraints are taken into account in the Gauss-Newton direction  $\xi_C(x)$  defined by (3.11), only those constraints in  $\hat{I}(x)$ , not aligned with the gradient  $\nabla J(x)$ , occur in the definition of  $\xi_J(x)$ .

*Remark 5.* Let us discuss two extreme cases related to the involved computational effort in the numerical implementation of (3.29). Upon discretization, we may assume that  $V = \mathbb{R}^k$  is a finite-dimensional space.

- (1) If the total number  $p + \tilde{q}$  of saturated or violated constraints is small compared to the dimension  $k$  of  $V$ , it is best, for numerical efficiency, to assemble the small square matrix  $(DC_{\hat{I}(x)} DC_{\hat{I}(x)}^\top)$  and to invert it by a direct method.
- (2) If  $V = \mathbb{R}^k$  is equipped with an inner product encoded by a matrix  $A$ , and if  $p + \tilde{q}$  is of the order of  $k$  or larger, the computation of the inverse of  $(DC_{\hat{I}(x)} DC_{\hat{I}(x)}^\top)$  can be expensive. However, it is still tractable if both  $DC$  and  $A$  are sparse matrices. For instance, this occurs in the case of bound constraints on the optimization variable  $x = (x_1, \dots, x_k)$ , e.g. constraints of the form  $\alpha_i \leq x_i \leq \beta_i$ ,  $i = 1, \dots, k$ . Recalling from Remark 1 that in this setting,  $DC_{\hat{I}(x)}^\top = A^{-1} DC_{\hat{I}(x)}^\top$ , it can be verified that the vector

$$X := A^{-1} DC_{\hat{I}(x)}^\top (DC_{\hat{I}(x)} A^{-1} DC_{\hat{I}(x)}^\top)^{-1} DC_{\hat{I}(x)} \nabla J(x)$$

can be computed as the solution of the sparse linear system

$$\begin{bmatrix} A & -DC_{\hat{I}(x)}^\top \\ DC_{\hat{I}(x)} & 0 \end{bmatrix} \begin{bmatrix} X \\ \Lambda \end{bmatrix} = \begin{bmatrix} 0 \\ DC_{\hat{I}(x)} \nabla J(x) \end{bmatrix},$$

where  $\Lambda \in \mathbb{R}^{p+\text{Card}(\tilde{I}(x))}$  is an extra slack variable, which yields the null space directions  $\xi_J(x) = \nabla J(x) - X$ . A similar strategy can be used to compute the range space direction  $\xi_C(x)$  of (3.11), or to solve the dual quadratic subproblem (3.13) by exploiting the sparsity of  $A$  and  $\text{DC}_{\tilde{I}(x)}$ .

*Remark 6.* As we have already mentioned, the Lagrange multiplier  $\mu^*(x)$  given by (3.21) may be understood as an indicator of which inequality constraints are aligned with the gradient of  $J$ . To further highlight this, it is instructive to consider the particular situation where the gradients of the constraint functions are orthogonal, i.e.:

$$\begin{aligned} a(\nabla g_i(x), \nabla g_j(x)) &= 0, \text{ for } i, j = 1, \dots, p, \ i \neq j, \\ a(\nabla h_i(x), \nabla h_j(x)) &= 0, \text{ for } i, j = 1, \dots, q, \ i \neq j, \\ a(\nabla g_i(x), \nabla h_j(x)) &= 0, \text{ for } i = 1, \dots, p, \ j = 1, \dots, q. \end{aligned}$$

Indeed, in this case, it easily follows from the Pythagore theorem that for any  $(\lambda, \mu) \in \mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x)}$ ,

$$\begin{aligned} \|\nabla J(x) + \text{Dg}(x)^T \lambda + \text{Dh}_{\tilde{I}(x)}(x)^T \mu\|_V^2 &= \left\| \nabla J(x) + \sum_{i=1}^p \lambda_i \nabla g_i(x) + \sum_{j \in \tilde{I}(x)} \mu_j \nabla h_j(x) \right\|_V^2 \\ &= \|\nabla J(x)\|_V^2 + \sum_{i=1}^p (\lambda_i^2 \|\nabla g_i(x)\|_V^2 + 2\lambda_i a(\nabla J(x), \nabla g_i(x))) + \sum_{j \in \tilde{I}(x)} (\mu_j^2 \|\nabla h_j(x)\|_V^2 + 2\mu_j a(\nabla J(x), \nabla h_j(x))). \end{aligned}$$

Therefore the minimization problem (3.13) is separable with respect to the variables  $(\lambda, \mu) \in \mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x)}$ :  $(\lambda_i^*(x))_{1 \leq i \leq p}$  and  $(\mu_i^*(x))_{i \in \tilde{I}(x)}$  are the respective solutions to the minimization problems:

$$\begin{aligned} \forall i \in 1 \dots p, \quad \lambda_i^*(x) &= \arg \min_{t \in \mathbb{R}} (t^2 \|\nabla g_i(x)\|_V^2 + 2ta(\nabla J(x), \nabla g_i(x))), \\ \forall i \in \tilde{I}(x), \quad \mu_i^*(x) &= \arg \min_{\substack{t \in \mathbb{R} \\ t \geq 0}} (t^2 \|\nabla h_i(x)\|_V^2 + 2ta(\nabla J(x), \nabla h_i(x))), \end{aligned}$$

which yields eventually:

$$\lambda_i^*(x) = -\frac{a(\nabla J(x), \nabla g_i(x))}{\|\nabla g_i(x)\|_V^2}, \quad \mu_i^*(x) = \begin{cases} 0 & \text{if } a(\nabla J(x), \nabla h_i(x)) \geq 0, \\ -\frac{a(\nabla J(x), \nabla h_i(x))}{\|\nabla h_i(x)\|_V^2} & \text{otherwise.} \end{cases}$$

Hence,  $\mu_i^*(x)$  is positive if and only if following the descent direction  $-\nabla J(x)$  leads to an increase (i.e. violation) of the  $i^{\text{th}}$  inequality constraint.

In the general case where all the constraint gradients are not mutually orthogonal, the interpretation of  $\mu^*(x)$  is similar, up to the additional complication that (3.13) accounts for the combinatorics behind the possible alignments between different constraint gradients. In the following, with a slight abuse of language, we shall nevertheless refer to the indices  $i \in \tilde{I}(x) \setminus \hat{I}(x)$  as those associated to constraints which are ‘aligned’ with  $\nabla J(x)$  (in the sense that  $-\text{Dh}_i(x)\xi_J(x) \leq 0$ , i.e. the violation  $h_i(x)$  decreases along  $-\xi_J(x)$  or, at worst, stay constant).

### 3.3.3. Behavior of the trajectories of the flow

The following proposition is the counterpart of Proposition 1 in the case of the equality and inequality constrained optimization problem (3.1).

**Proposition 5.** Assume the trajectories  $x(t)$  of the flow

$$\begin{cases} \dot{x}(t) = -\alpha_J \xi_J(x(t)) - \alpha_C \xi_C(x(t)) \\ x(0) = x_0, \end{cases} \quad (3.30)$$

with  $\xi_J$  and  $\xi_C$  given by (3.11) and (3.29) exist on some interval  $[0, T]$  for  $T > 0$  and are such that:

(a) The set  $\tilde{I}(x(t))$  defined in (3.2) is constant over  $[0, T]$ :

$$\forall t \in [0, T], \quad \tilde{I}(x(t)) = \tilde{I}(x_0)$$

(b) The constraints remain qualified along the flow  $x(t)$ , in the sense of (3.3).

Then the following properties hold true:

(1) The violation of the constraints decreases exponentially:

$$\forall t \in [0, T], \mathbf{g}(x(t)) = e^{-\alpha_C t} \mathbf{g}(x_0) \text{ and } \mathbf{h}(x(t)) \leq e^{-\alpha_C t} \mathbf{h}(x_0). \quad (3.31)$$

(2) Assume that  $\text{rank}(\text{DC}_{\tilde{I}(x_0)}(x))$  is maximal for all  $x$  in  $K = \{x \in V \mid \|\mathbf{C}_{\tilde{I}(x_0)}(x)\|_\infty \leq \|\mathbf{C}_{\tilde{I}(x_0)}(x_0)\|_\infty\}$  and

$$\sup_{x \in K} \|\nabla J(x)\|_V |\sigma_p^{-1}(x)| < +\infty. \quad (3.32)$$

where  $\sigma_p(x)$  is the smallest singular value of  $\text{DC}_{\tilde{I}(x)}(x)$ . Then  $J$  decreases as soon as the projected gradient is large with respect to the violation of the constraints (which goes to 0 according to (3.31)), namely there exists a constant  $C > 0$  such that

$$\forall t \in [0, T], \|\Pi_{\mathbf{C}_{\tilde{I}(x(t))}}(\nabla J(x(t)))\|_V^2 > C e^{-\alpha_C t} \Rightarrow \frac{d}{dt} J(x(t)) < 0. \quad (3.33)$$

(3) Any stationary point  $x^*$  of the flow (3.30) satisfies the KKT optimality conditions (3.6) which equivalently rewrite:

$$\begin{cases} \nabla J(x^*) + \text{D}\mathbf{g}(x^*)^\top \boldsymbol{\lambda}^*(x^*) + \text{D}\mathbf{h}_{\tilde{I}(x^*)}(x^*)^\top \boldsymbol{\mu}^*(x^*) = 0, \\ \mathbf{g}(x^*) = 0 \text{ and } \mathbf{h}_{\tilde{I}(x^*)}(x^*) = 0 \Leftrightarrow \mathbf{C}_{\tilde{I}(x^*)}(x^*) = 0, \end{cases} \quad (3.34)$$

where  $(\boldsymbol{\lambda}^*(x^*), \boldsymbol{\mu}^*(x^*)) \in \mathbb{R}^p \times \mathbb{R}_+^{\tilde{q}(x^*)}$  are defined in (3.13) or (3.21).

*Proof.*

(1) Definition (3.11) of  $\boldsymbol{\xi}_C(x(t))$  implies  $\text{DC}_{\tilde{I}(x(t))} \boldsymbol{\xi}_C(x(t)) = \mathbf{C}_{\tilde{I}(x(t))}(x(t))$ , and since  $-\boldsymbol{\xi}_J(x(t))$  is positively proportional to  $\boldsymbol{\xi}^*(x(t))$  (Proposition 3), it holds

$$\text{DC}_{\tilde{I}(x(t))} \boldsymbol{\xi}_J(x(t)) = 0, \quad -\text{D}\mathbf{h}_{\tilde{I}(x(t)) \setminus \hat{I}(x(t))}(x(t)) \boldsymbol{\xi}_J(x(t)) \leq 0.$$

Therefore we obtain

$$\frac{d}{dt} \mathbf{C}_{\hat{I}(x(t))}(x(t)) = -\alpha_C \mathbf{C}_{\hat{I}(x(t))}(x(t)) \text{ and } \frac{d}{dt} \mathbf{h}_{\tilde{I}(x(t)) \setminus \hat{I}(x(t))}(x(t)) \leq -\alpha_C \mathbf{h}_{\tilde{I}(x_0) \setminus \hat{I}(x_0)}(x(t)) \quad (3.35)$$

from which (3.31) follows by application of Gronwall's lemma.

(2) The proof is identical to that of Proposition 1.

(3) A stationary point  $x^*$  of (3.30) satisfies by definition  $-\alpha_J \boldsymbol{\xi}_J(x^*) - \alpha_C \boldsymbol{\xi}_C(x^*) = 0$ . Left multiplication of this identity by  $\text{DC}_{\tilde{I}(x^*)}(x^*)$  yields:

$$-\alpha_J \text{DC}_{\tilde{I}(x^*)}(x^*) \boldsymbol{\xi}_J(x^*) - \alpha_C \mathbf{C}_{\tilde{I}(x^*)}(x^*) = 0. \quad (3.36)$$

Remembering now that from definition (3.12),

$$-\text{DC}_{\tilde{I}(x^*)} \boldsymbol{\xi}_J(x^*) \leq 0 \text{ and } \mathbf{C}_{\tilde{I}(x^*)}(x^*) \geq 0,$$

equality in (3.36) can hold only if both terms vanish. In particular, we infer that  $\mathbf{C}_{\tilde{I}(x^*)}(x^*) = 0$ , a fact which implies  $\boldsymbol{\xi}_C(x^*) = 0$  and which encompasses the last two lines of the KKT conditions (3.6). Returning to the fact that  $-\alpha_J \boldsymbol{\xi}_J(x^*) - \alpha_C \boldsymbol{\xi}_C(x^*) = 0$ , we obtain that  $\boldsymbol{\xi}_J(x^*) = 0$ , which corresponds to the first line in (3.6). This completes the proof.  $\square$

*Remark 7.* The assumption (a) in Proposition 5, whereby the index set  $\tilde{I}(x(t))$  remains constant is essentially made to ensure that the right-hand side of the flow (3.30) is continuous. Indeed, in such a case, the range space direction  $\boldsymbol{\xi}_C(x(t))$  is continuous by its definition (3.11), while the null space step  $\boldsymbol{\xi}_J(x(t))$  is continuous because

$$\boldsymbol{\xi}_J(x(t)) = \nabla J(x(t)) + \text{DC}_{\tilde{I}(x(t))} \begin{bmatrix} \boldsymbol{\lambda}^*(x(t)) \\ \boldsymbol{\mu}^*(x(t)) \end{bmatrix}$$

and it can be shown that the multipliers  $(\boldsymbol{\lambda}^*(x(t)), \boldsymbol{\mu}^*(x(t)))$  defined by (3.13) are continuous functions. At a time  $T$  corresponding to a sudden change of the index set  $\tilde{I}(x(t))$ , we assume that the solution  $x(t)$  can

be extended by restarting the ODE (3.30) with the new index set  $\tilde{I}(x(T))$ . From (1) in Proposition 5, the bound  $\mathbf{h}(x(t)) \leq e^{-\alpha_C t} \mathbf{h}(x(0))$  still holds after the time  $T$  for all constraints  $i \in \{1, \dots, q\}$ : constraints are asymptotically satisfied. Properties (2) and (3) remain true, up to an adjustment of the constant  $C$  in (3.33) (which can be taken global since there are finitely many possible sets  $\tilde{I}(x(t))$ ). There may exist situations where the set of asymptotically saturated constraints  $\tilde{I}(x(t))$  could oscillate indefinitely. However (2) states that  $x(t)$  always keeps improving (in the sense of (3.33)), and (3) states that if  $x(t)$  eventually converges, it is necessarily towards a KKT point.

*Remark 8.* In practice, the analysis of Proposition 5 is sufficient because, similarly to the conclusions of Remark 2, analogous properties hold for the discrete scheme

$$x_{n+1} = x_n - \Delta t (\alpha_J \boldsymbol{\xi}_J(x_n) + \alpha_C \boldsymbol{\xi}_C(x_n)). \quad (3.37)$$

Indeed, one can easily check that:

- (1) Up to first order, the violation of the constraints decreases at a geometric rate:

$$\mathbf{C}(x_{n+1}) = (1 - \alpha_C \Delta t) \mathbf{C}(x_n) + o(\Delta t). \quad (3.38)$$

This suggests that in order to obtain a stable scheme, one must a priori select  $\alpha_C$  and  $\Delta t$  such that  $0 < \alpha_C \Delta t < 2$ .

- (2)  $x^*$  is an accumulation point of the sequence  $(x_n)_{n \in \mathbb{N}}$  if and only if it is feasible, i.e.  $\mathbf{C}_{\tilde{I}(x^*)}(x^*) = 0$  and  $x^*$  is a KKT point.

Finally, note that a flexibility of this ODE approach is that at the continuous level, the results of Proposition 5 do not depend on the values of the parameters  $\alpha_J > 0$  and  $\alpha_C > 0$ . Therefore the convergence of the discrete scheme towards the continuous trajectory should hold as soon as the discretization step size  $\Delta t > 0$  is sufficiently small.

### 3.4. Comparison between the proposed method and the use of slack variables

The main differences between the slack variable approach of Section 3.2 and the proposed flow (3.30) in Section 3.3 for dealing with equality and inequality constrained problems can be summarized as follows.

- (1) Any point  $x^{\text{crit}}$  satisfying the constraints ( $\mathbf{C}_{\tilde{I}(x^{\text{crit}})}(x^{\text{crit}}) = 0$ ) and  $\Pi_{\mathbf{C}_{\tilde{I}(x^{\text{crit}})}}(\nabla J(x^{\text{crit}})) = 0$  is a stationary point of the extended dynamical system (3.8), although it might violate the full KKT condition (because (3.5) may yield possible negative values of the multiplier  $\boldsymbol{\mu}_{\tilde{I}(x^{\text{crit}})}(x^{\text{crit}})$ ). In contrast,  $x^*$  is a stationary points of the flow (3.30) if and only if it is a true feasible KKT point, see Proposition 5.
- (2) The computation of  $\boldsymbol{\xi}_J(x)$  and  $\boldsymbol{\xi}_C(x)$  in our flow (3.30) requires to invert a matrix of size at most  $(p + \tilde{q}(x))$ -by- $(p + \tilde{q}(x))$  with  $\tilde{q}(x)$  the number of active or violated constraints at  $x$ . The process of equalizing inequality constraints as in [38, 41] rather requires to invert the full  $(p + q)$ -by- $(p + q)$  matrix  $\mathbf{D}\mathbf{C}(x, z)\mathbf{D}\mathbf{C}(x, z)^T$ . Our method is therefore more efficient if  $\tilde{q}(x) \ll q$ , that is if a lot of inequality constraints are inactive.
- (3) At feasible points, our ODE (3.30) follows the best locally admissible descent direction (with respect to the norm of  $V$ ). This is not the case for the extended ODE (3.8). Therefore, from a common feasible point  $x$ , the ODE (3.30) always decreases the objective function with a steeper slope  $dJ/ds$  with respect to the parameterization induced by the path length  $s$ , defined as a function of the time  $t$  by

$$s(t) = \int_0^t \|\dot{x}(\alpha)\|_V d\alpha. \quad (3.39)$$

This property is illustrated in the academic examples of Section 5, and in particular on Figure 8a below.

All in all, our observations based on the simple numerical examples of the next section tend to illustrate that both flows (3.8) and (3.30) may have equivalent performances for solving the non linear optimization problem (1.1), this performance being measured in term of the total length

$$S = \int_0^{+\infty} \|\dot{x}\|_V dt$$

20

covered by the optimization path to reach the optimum. However, the two ODEs (3.8) and (3.30) yield optimization paths of essentially different natures. Our null space flow (3.30) ignores inactive constraints and those aligned with the gradient of the objective function. As a result, it produces non smooth paths that are more likely to reach quickly the saturation of the constraint. The extended flow (3.8) yields smoother trajectories that more likely stay away from the constraints, at the cost of inverting at every step the full matrix  $\mathbf{DC}(x, z)\mathbf{DC}(x, z)^T$  of the size of the total number of constraints (active and inactive).

#### 4. PRACTICAL IMPLEMENTATION DETAILS

In this section, the ODE (1.3) is discretized by an explicit scheme and we propose a generic algorithm for adapting the time step  $\Delta t$ . We also discuss small adaptations for accounting for discontinuous changes of the right-hand side  $-(\alpha_J \xi_J + \alpha_C \xi_C)$ .

##### 4.1. Time step adaptation based on a merit function.

The ODE (1.3) is discretized by an explicit scheme of the form:

$$x_{n+1} = x_n - \Delta t_n (\alpha_J \xi_J(x_n) + \alpha_C \xi_C(x_n)), \quad (4.1)$$

with a variable time step  $\Delta t_n > 0$ . The practical implementation of such a strategy is often guided by a merit function, i.e. an indicator allowing to detect that a step has been too large, a situation where a choice has to be made regarding whether the step should be reduced or accepted. For our null space algorithm, a merit function which resembles very much that of the Augmented Lagrangian Method is readily available, however with a specific choice of multipliers:

**Lemma 3.** *For a given  $x_n \in V$ , let  $\text{merit}_{x_n} : V \rightarrow \mathbb{R}$  be the function defined by*

$$\text{merit}_{x_n}(x) := \alpha_J \left( J(x) + \mathbf{\Lambda}(x_n)^T \mathbf{C}_{\tilde{I}(x_n)}(x) \right) + \frac{\alpha_C}{2} \mathbf{C}_{\tilde{I}(x_n)}(x)^T \mathbf{S}(x_n) \mathbf{C}_{\tilde{I}(x_n)}(x) \quad (4.2)$$

where  $\mathbf{\Lambda}(x_n) = \begin{bmatrix} \lambda^*(x_n)^T & \mu^*(x_n)^T \end{bmatrix}^T$  is the vector of multipliers defined as the solution to the dual problem (3.13) (see (3.21)) and  $\mathbf{S}(x_n) = (\mathbf{DC}_{\tilde{I}(x_n)}(x_n) \mathbf{DC}_{\tilde{I}(x_n)}(x_n)^T)^{-1}$  is symmetric positive definite. Then (4.1) is a gradient step for decreasing the function  $\text{merit}_{x_n}$ , namely:

$$\nabla \text{merit}_{x_n}(x_n) = \alpha_J \xi_J(x_n) + \alpha_C \xi_C(x_n).$$

*Proof.* It is a straightforward computation of the gradient of (4.2). □

A possible implementation of an optimization strategy of the form (4.1) based on this merit function is summarized in Algorithm 1, which requires the introduction of a few extra parameters:

- **time\_step**: choose a fixed time step  $\Delta t > 0$ .
- **maxtrials**: the optimization time step is decreased up to **maxtrials** times until the value of the merit function has decreased. If the merit function has not decreased after all **maxtrials** steps, the smallest step is accepted.
- **tolLag**: a small threshold for the values of the Lagrange multipliers  $\mu_i^*$  under which these are considered to be 0 (we took **tolLag**=1e-8). This value should be set in accordance with the machine precision and that of the quadratic programming solver for the dual problem (3.13).

Let us emphasize that these parameters have a quite intuitive and physical meaning, so that the task of assigning their values does not involve fine tunings in practice.

Importantly, the rescaling induced by the inverse of the correlation matrix  $(\mathbf{DC}_{\tilde{I}(x_n)} \mathbf{DC}_{\tilde{I}(x_n)}^T)^{-1}$  normalizes all the constraints; in particular, the whole Algorithm 1 is invariant under multiplication of the constraints by arbitrary positive constants (up to the machine precision for the step 3); a preliminary rescaling of the constraints is therefore not required from the user.



---

**Algorithm 1** Discretization of the null space gradient flow (3.30), based on a merit function.

---

```

1: for  $n = 1 \dots \text{maxiter}$  do
2:   Determine  $\tilde{I}(x_n) = \{i \in \{1, \dots, q\} \mid h_i(x_n) \geq 0\}$  the set of active or violated constraints.
3:   Solve the dual problem (3.13) to obtain the Lagrange multiplier  $\mu^*(x_n)$ . Infer the subset  $\hat{I}(x_n) \subset \tilde{I}(x_n)$  of Proposition 4, indicating which constraints must be active in (3.12):
      
$$\hat{I}(x_n) = \{i \in \tilde{I}(x_n) \mid \mu_i^*(x_n) > \text{tolLag}\}. \quad (4.3)$$

4:   Extract the vectors  $C_{\hat{I}(x_n)}(x_n)$  and  $C_{\tilde{I}(x_n)}(x_n)$  (defined by (1.6)) and compute
      
$$\begin{aligned} \xi_J(x_n) &= (I - DC_{\hat{I}(x_n)}^T)(DC_{\hat{I}(x_n)}DC_{\hat{I}(x_n)}^T)^{-1}DC_{\hat{I}(x_n)}\nabla J(x_n), \\ \xi_C(x_n) &= DC_{\tilde{I}(x_n)}^T(DC_{\tilde{I}(x_n)}DC_{\tilde{I}(x_n)}^T)^{-1}C_{\tilde{I}(x_n)}. \end{aligned} \quad (4.4)$$

5:   for  $k = 1 \dots \text{maxtrials}$  do
6:     Perform the step
      
$$x_{n+1} = x_n - \frac{\Delta t}{2^{k-1}}(\alpha_J \xi_J(x_n) + \alpha_C \xi_C(x_n)).$$

7:     if  $\text{merit}_{x_n}(x_{n+1}) < \text{merit}_{x_n}(x_n)$  then
8:       break
9:     end if
10:  end for
11: end for

```

---

#### 4.2. Accounting for discontinuities near the inequality constraint barriers

A potential issue when implementing the above Algorithm 1 comes from the fact that the vector fields  $\xi_J$  and  $\xi_C$  given by (3.11) and (3.29) are characterized by the same discontinuities as the discrete index mapping  $x \mapsto \tilde{I}(x)$ . As a result, abrupt oscillations of the discrete optimization path  $(x_n)$  may occur near the boundary of the feasible set: if  $h_i(x_n) = 0$  and  $i \in \tilde{I}(x_n)$  for some index  $i \in \{1, \dots, q\}$ , then in the definition (3.29) of  $\xi_J(x_n)$ , the gradient  $\nabla J(x_n)$  is projected tangentially to the constraint  $h_i$ , but it is not projected after any slight deviation (e.g. due to the discretization) making this constraint inactive ( $h_i(x_{n+1}) < 0$ ). This kind of issue is very classical in the discretization of ODEs with discontinuous vector fields and can be tackled by various methods, see e.g. [21] for a review.

In this section, we suggest a simple alternative: constraints are felt from a short distance by replacing the set  $\tilde{I}(x_n)$  in (1.9) (step 3 of Algorithm 1) with the set  $\tilde{I}_\epsilon(x_n)$  of inequality constraints violated “up to  $\epsilon_i$ ”:

$$\tilde{I}_\epsilon(x_n) = \{i \in \{1, \dots, q\} \mid h_i(x_n) \geq -\epsilon_i\}. \quad (4.5)$$

The tolerances  $\epsilon_i > 0$  can be estimated in an automatic fashion, independent of an arbitrary rescaling of the constraints, thanks to an a posteriori bound we now detail. Let  $\mathbf{h}$  be a user-defined parameter accounting for the distance from the optimization path at which the constraints should be felt. This characteristic length  $\mathbf{h}$  should be defined in accordance with the typical distance  $\|\Delta x\|_V = \|x_{n+1} - x_n\|_V$  between two successive iterations; in the academic examples in Section 5 below considering optimization in  $\mathbb{R}^k$ , we may set e.g.  $\mathbf{h} = 0.01$  for a typical increment size  $\|\Delta x\|_V \simeq 0.1$ . For our shape optimization applications in Section 6,  $\mathbf{h}$  is typically of the order of the discretization mesh size, see Section 6.2.2 below.

Assume now that the current point  $x_n$  satisfies the constraint  $h_i$  up to the uncertainty  $\mathbf{h}$  on its location: by this we mean that there exists some unknown point  $x_n^*$  such that  $\|x_n^* - x_n\| \leq \mathbf{h}$ ,  $h_i(x_n) > 0$  and  $h_i(x_n^*) = 0$ . Then the error  $\mathbf{h}$  for the location of  $x_n$  propagates to the constraint values  $h_i(x_n)$  according to the following inequality:

$$h_i(x_n) = |h_i(x_n) - h_i(x_n^*)| \simeq |Dh_i(x_n)(x_n^* - x_n)| \leq \|\nabla h_i(x_n)\|_V \mathbf{h}. \quad (4.6)$$

It is therefore natural to set

$$\epsilon_i := \|\nabla h_i(x_n)\|_V \mathbf{h} \quad (4.7)$$

for the value of  $\epsilon_i$  in (4.5). Note that more generally, the a posteriori bound (4.6) allows to assert whether a constraint  $C_i(x_n)$  can be considered as satisfied or not with respect to the numerical discretization.

The dual problem (3.13) is then solved with  $\tilde{I}_\epsilon(x_n)$  instead of  $\tilde{I}(x_n)$  in order to obtain a new subset of indices  $\hat{I}_\epsilon(x_n)$  which indicates which constraints are likely to be not aligned with the gradient  $\nabla J(x_n)$  when crossing the barrier  $\mathbf{h} = 0$ . The null space and range space steps  $\xi_J(x_n)$  and  $\xi_C(x_n)$  in step 4 of Algorithm 1 are finally replaced with  $\xi_{J,\epsilon}(x_n)$  and  $\xi_{C,\epsilon}(x_n)$  computed as follows:

$$\xi_{J,\epsilon}(x_n) := (\mathbf{I} - \text{DC}_{\hat{I}_\epsilon(x_n)}^T)(\text{DC}_{\hat{I}_\epsilon(x_n)}\text{DC}_{\hat{I}_\epsilon(x_n)}^T)^{-1}\text{DC}_{\hat{I}_\epsilon(x_n)}\nabla J(x_n), \quad (4.8)$$

$$\xi_{C,\epsilon} := \text{DC}_{I_\epsilon^*(x_n)}^T(\text{DC}_{I_\epsilon^*(x_n)}\text{DC}_{I_\epsilon^*(x_n)}^T)^{-1}\mathbf{C}_{I_\epsilon^*(x_n)}(x_n), \quad (4.9)$$

where  $I_\epsilon^*(x_n) = \tilde{I}(x_n) \cup \hat{I}_\epsilon(x_n)$  is the set of constraints that are either violated, saturated or not aligned with the gradient  $\nabla J(x_n)$  at  $\mathbf{h} = -(\epsilon_1, \dots, \epsilon_q)^T$ . The use of  $\hat{I}_\epsilon(x_n)$  in the definition of  $\xi_{J,\epsilon}(x_n)$  ensures that the gradient  $\nabla J(x_n)$  is being projected tangentially to the constraint on a small layer near the boundary of the feasible set. As a result, no abrupt discontinuity occurs anymore for  $\xi_{J,\epsilon}$  and  $\xi_{C,\epsilon}$  when crossing the boundary of the feasible domain while remaining in this layer. Including constraints  $i \in \hat{I}_\epsilon(x_n)$  in the Gauss-Newton direction  $\xi_{C,\epsilon}(x_n)$  even if they are satisfied (i.e. if  $-\epsilon_i \leq h_i(x_n) \leq 0$ ) further allows to stabilize the values of these constraints closer to zero.

## 5. ILLUSTRATIONS AND COMPARISONS ON ACADEMIC TEST CASES

In this section, we consider simple and illustrative academic examples in order to compare qualitatively the following three strategies for dealing with inequality constraints in optimization problems:

- The method of Section 3.2 for equalizing inequality constraints by means of slack variables, see (3.30); this strategy is hereafter labelled as ‘SLACK’.
- The proposed null space flow (3.30) in Section 3.3, based on the dual problem (3.13) for solving the combinatorial character of the constraints, which is labelled as ‘NLSPACE’.
- An alternative, naive version of (3.30) which does not take advantage of the use of a dual problem, and simply projects  $\nabla J(x(t))$  on all the violated constraints:

$$\begin{cases} \dot{x} = -\alpha_J \tilde{\xi}_J(x(t)) - \alpha_C \xi_C(x(t)) \\ \tilde{\xi}_J(x) := (\mathbf{I} - \text{DC}_{\tilde{I}(x)}^T)(\text{DC}_{\tilde{I}(x)}\text{DC}_{\tilde{I}(x)}^T)^{-1}\text{DC}_{\tilde{I}(x)}\nabla J(x) \\ \xi_C(x) := \text{DC}_{\tilde{I}(x)}^T(\text{DC}_{\tilde{I}(x)}\text{DC}_{\tilde{I}(x)}^T)^{-1}\mathbf{C}_{\tilde{I}(x)}(x). \end{cases} \quad (5.1)$$

In other words, all the violated or saturated constraints are taken into account in the computation of both the null space and range space directions  $\xi_J(x)$  and  $\xi_C(x)$ . This strategy is labelled as ‘NLSPACE (no dual)’.

To achieve our comparison purpose, Algorithm 1 is implemented for the discretization of (3.8), (3.30) and (5.1), with straightforward adaptations for equalizing slack variables or disabling the resolution of the dual problem. In all the following cases considered, we have set the values of  $\alpha_J$  and  $\alpha_C$  such that  $\alpha_J/\alpha_C = 5/3$ . The step size  $\Delta t$  was chosen sufficiently small to compute continuous paths with satisfying accuracy. Our discussion is exclusively focused on the continuous trajectories of the considered ODEs. In particular, we do neither discuss the issue of the selection of the time step, nor the efficiency of these methods in terms of the needed number of iterations required to achieve convergence.

In order to compare the three methods without bias, we consider the arc length  $s(t)$  (defined in (3.39)) as the common reference time for the three ODEs (3.8), (3.30) and (5.1); recall indeed that this quantity is invariant under any monotone parameterization change of the time  $t$ . In the convergence figures below, optimized quantities are then plotted with respect to the pseudo time  $s(t)$  in abscissa, for example we plot the graph  $t \mapsto (s(t), J(s(t)))$  in order to account for the evolution of the objective function  $J$ .

We shall also plot the evolution of the Lagrange multipliers  $s \mapsto \mu(x(s))$  associated with  $\xi_J(x(s))$  or  $\tilde{\xi}_J(x(s))$  for the ODEs (3.30) and (5.1). For that purpose, these Lagrange multipliers are defined on the violated indices  $i \in \tilde{I}(x(s))$  by (3.21) for the null space flow (3.30), and by

$$\mu(x(s)) := -(\text{DC}_{\tilde{I}(x(s))}^T)^{-1}\text{DC}_{\tilde{I}(x(s))}\nabla J(x(s)) \quad (5.2)$$

for the flow (5.1) that does not use the dual problem (3.13). For the indices  $i \in \{1, \dots, q\} \setminus \tilde{I}(x(s))$ , the value of the Lagrange multiplier is set to  $\mu_i(x(s)) := 0$  by convention. We do not plot Lagrange multipliers for the ODE (3.8) using slack variables because these are defined with respect to the extended variables  $(x(s), z(s))$ .

The examples of this section take place in the optimization set  $\mathcal{X} = \mathbb{R}^2$ , which is equipped with the usual euclidean inner product; the Hilbert transposition  $\mathcal{T} = T$  coincides with the usual transposition operator (see Definition 1). For simplicity, these examples only involve inequality constraints; we consider the following three scenarios:

- In Section 5.1, the initial point is unfeasible and the gradient of the objective function  $\nabla J(x)$  is always aligned with the directions of the constraints;
- In Section 5.2 also, the initial point is unfeasible, but the gradient  $\nabla J(x)$  may not be aligned with the direction of constraints;
- In Section 5.3, one of the constraints becomes inactive in the course of the optimization path.

Further comparisons with other iterative optimization algorithms will be presented in the PhD thesis [26].

### 5.1. Test case 1 : unfeasible initialization with initial gradient aligned with the constraints.

Our first example features the following problem, reproduced from [25]:

$$\begin{aligned} \min_{(x_1, x_2) \in \mathbb{R}^2} \quad & J(x_1, x_2) := x_2 + 0.3x_1 \\ \text{s.t.} \quad & \begin{cases} h_1(x_1, x_2) := -x_2 + \frac{1}{x_1} \leq 0, \\ h_2(x_1, x_2) := x_1 + x_2 - 3 \leq 0. \end{cases} \end{aligned} \quad (5.3)$$

This test case is designed so that for the chosen initial point  $x_0 = (1.5, 2.25)$ , the gradient of the objective function  $\nabla J(x)$  is ‘aligned’ with the linear constraint  $h_2$ , in the sense that

$$-\nabla h_2(x_0) \cdot \nabla J(x_0) < 0.$$

Hence at least for small times (in fact during the whole optimization path), the constraint  $h_2$  can be ignored since the minimization of  $J$  is naturally concurrent with a decrease of the value of  $h_2$ .

The optimization paths taken by the solutions of the three ODEs (3.8), (3.30) and (5.1) are plotted on Figure 1. The associated convergence histories for the values of the objective and constraint functions are displayed on Figure 2.

Let us comment the trajectory followed by the null space flow (3.30) in details. The gradient of the objective function remains aligned with the constraint  $h_2$ , which is associated with a zero Lagrange multiplier  $\mu_2(x(s))$  (see Figure 3). During the first part of the optimization, the first constraint  $h_1$  is not violated, hence the multiplier  $\mu_1(x(s))$  is also set to zero. As a result, both constraints are ignored when computing the null space direction, which is set equal to the gradient:  $\xi_J(x(s)) = \nabla J(x(s))$ . The optimization path  $x(s)$  follows then almost the direction of the gradient  $\nabla J(x(s))$  (without projection), up to a small deviation induced by the non zero Gauss-Newton direction  $\xi_C(x(s))$  in the unfeasible domain. When the hyperbolic constraint represented by  $h_1$  becomes violated, the gradient is not aligned anymore with this constraint and the dual problem (3.13) yields a non-zero Lagrange multiplier  $\mu_1(x(s))$  (near  $s = 1.4$ ). From this point, the gradient  $\nabla J(x(s))$  is then projected tangentially to the constraint  $h_1$  till the optimum is attained.

In contrast, the path of the ODE (5.1) (no dual problem) fails to find the optimum as it is unable to unstick from the first saturated constraint. Notably, the gradient  $\nabla J(x)$  is kept being projected tangentially to the violated constraint  $h_2$  while it should not, which could have been detected from the negativity of the computed Lagrange multiplier  $\mu_2$  (see Figure 3).

Finally, the extended ODE (3.8) making use of slack variables feels the constraint  $h_1$  from distance, inducing a deviation of the trajectory in the direction of the optimum before reaching the saturation of the constraint. This allows the trajectory  $x(s)$  to find globally a slightly shorter path than our null space flow (3.30).

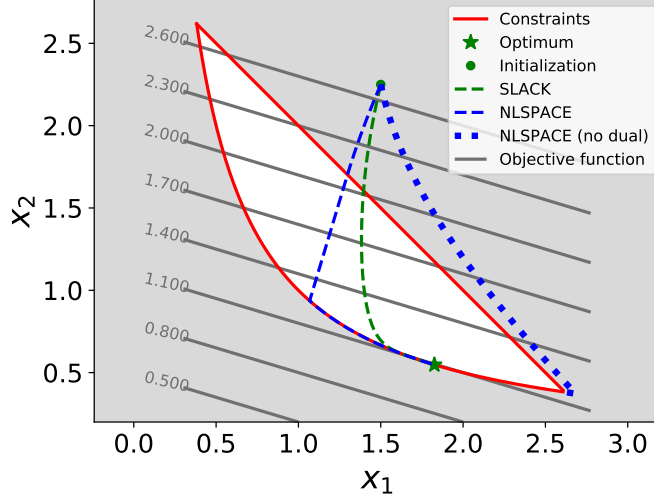
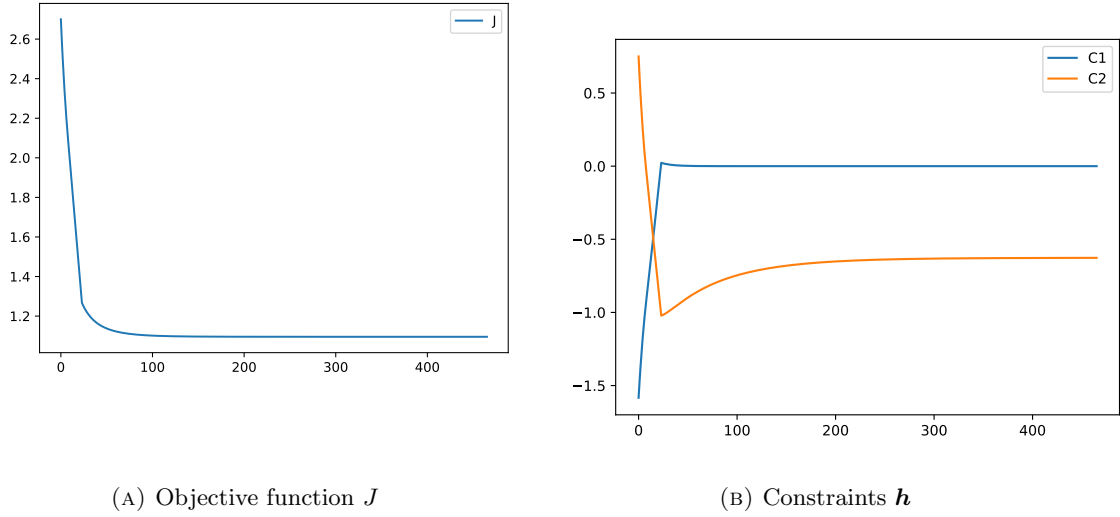


FIGURE 1. Optimization problem of Section 5.1: optimization paths for an unfeasible initialization  $x_0$  with  $\nabla J(x_0)$  aligned in the direction of the constraints.



(A) Objective function  $J$

(B) Constraints  $h$

FIGURE 2. History curves for the optimization problem of Section 5.1.

## 5.2. Test case 2 : unfeasible initialization with initial gradient not aligned with the constraints.

We now devise a test case where the gradient of the initialization is not aligned with the constraints. The feasible domain is the same as in the previous test case but the objective function is different:

$$\begin{aligned} \min_{(x_1, x_2) \in \mathbb{R}^2} \quad & J(x_1, x_2) := (x_1 - 2)^2 + (x_2 - 2)^2 \\ \text{s.t.} \quad & \begin{cases} h_1(x_1, x_2) := -x_2 + \frac{1}{x_1} \leq 0 \\ h_2(x_1, x_2) := x_1 + x_2 - 3 \leq 0. \end{cases} \end{aligned} \quad (5.4)$$

We keep the same initialization  $x_0 = (1.5, 2.25)$ . Corresponding optimization paths and convergence curves are displayed on Figs. 4 and 5.

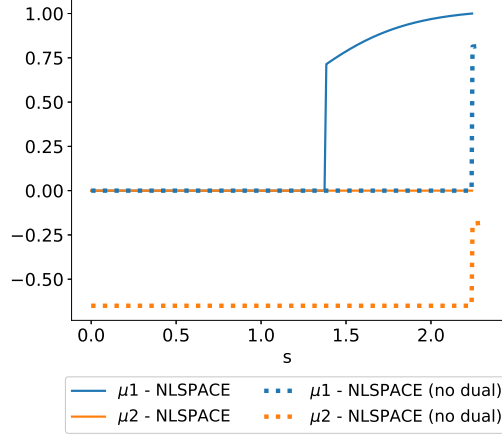


FIGURE 3. Evolution of the Lagrange multipliers  $\mu_1(x(s)), \mu_2(x(s))$  for the optimization test case of Section 5.1.

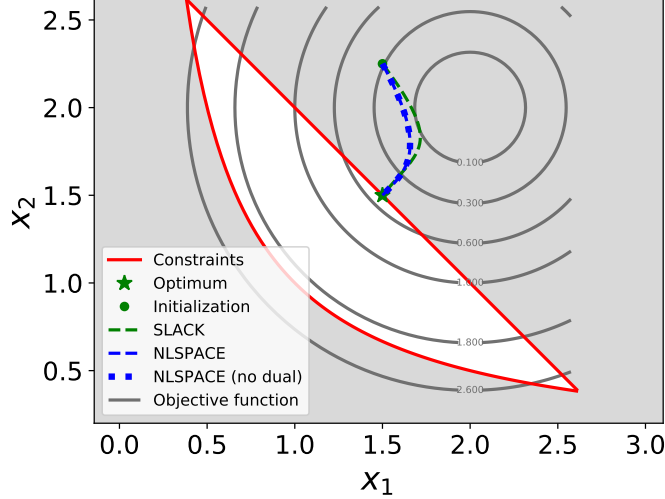


FIGURE 4. Optimization problem of Section 5.2: unfeasible initialization  $x_0$  with  $\nabla J(x_0)$  not aligned in the direction of the constraints.

For this example, the linear constraint  $h_2$  is not aligned with the gradient along the optimization path of the null space gradient flow (3.30). This is associated with a non-zero Lagrange multiplier  $\mu_2(x(s)) > 0$  (see Figure 6): the gradient  $\nabla J(x(s))$  is kept being projected tangentially to the constraint  $h_2$  when computing  $\xi_J(x(s))$ . For this case, the combination with the Gauss-Newton direction  $\xi_C(x(s))$  allows to decrease simultaneously the objective function and the violation of the constraints, which enables the optimization path to reach directly the optimum when hitting the feasible set. Note that the convergence curve Figure 5a depicts a monotonically increasing objective function  $J$  throughout the optimization path (although we are minimizing  $J$ ); this is of course due to the fact that constraints are never satisfied.

Since the constraint  $h_2$  remains active from initialization to the optimum, the optimization path is unchanged when disabling the dual problem (ODE (5.1)). Finally, the path selected by the extended flow (3.8) to reach the optimum is very similar to the one of our method, although slightly longer.

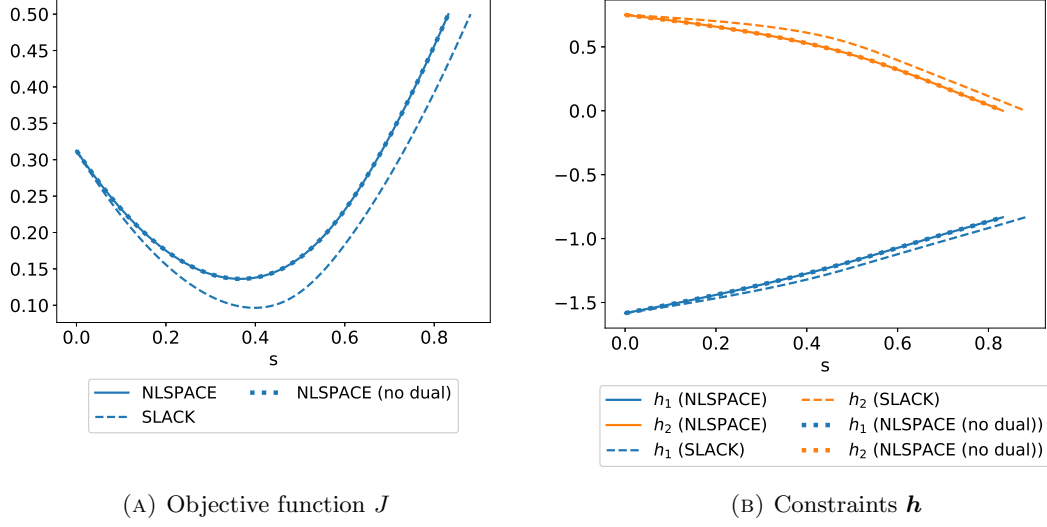


FIGURE 5. History curves for the optimization problem of Section 5.2.

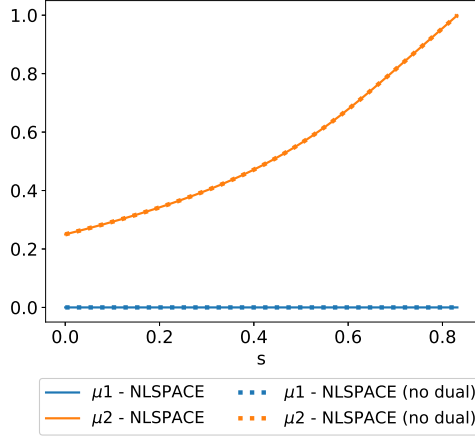


FIGURE 6. Evolution of the Lagrange multipliers  $\mu_1(x(s)), \mu_2(x(s))$  for the optimization test case of Section 5.2.

### 5.3. Test case 3: a saturated inequality constraint becoming inactive along the optimization path

This last optimization test case is designed to illustrate the usefulness of the dual problem for detecting when a saturated inequality constraint becomes unsaturated. We consider a disconnected unfeasible domain made from the reunion of a half-space and the interior region of a parabola:

$$\begin{aligned} \min_{(x_1, x_2) \in \mathbb{R}^2} \quad & J(x_1, x_2) = x_1^2 + (x_2 + 3)^2 \\ \text{s.t.} \quad & \begin{cases} h_1(x_1, x_2) = -x_1^2 + x_2 \leq 0 \\ h_2(x_1, x_2) = -x_1 - x_2 - 2 \leq 0 \end{cases} \end{aligned} \quad (5.5)$$

The feasible domain and optimization paths starting from the initialization  $x_0 = (3, 3)$  are displayed on Figure 7. Associated convergence curves are reported on Figure 8.

For the null space flow (3.30), four different stages occur as is visible on the evolution of the Lagrange multipliers  $\mu_1(x(s)), \mu_2(x(s))$  reported on Figure 9:

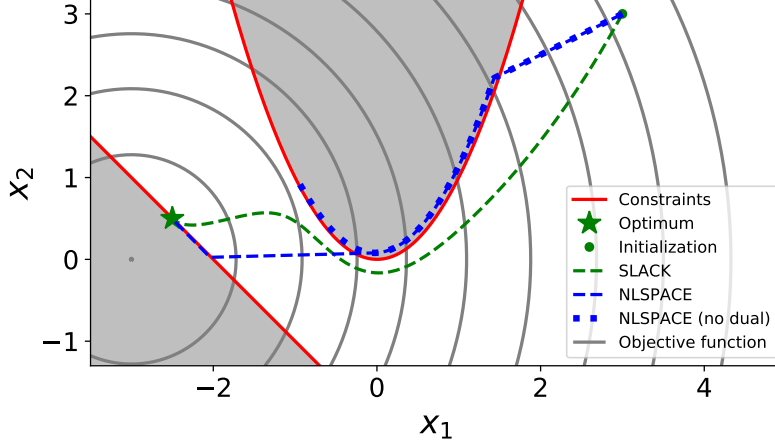


FIGURE 7. Optimization problem of Section 5.3 : feasible initialization  $x_0$  but the optimization has to find a path across the parabolic domain.

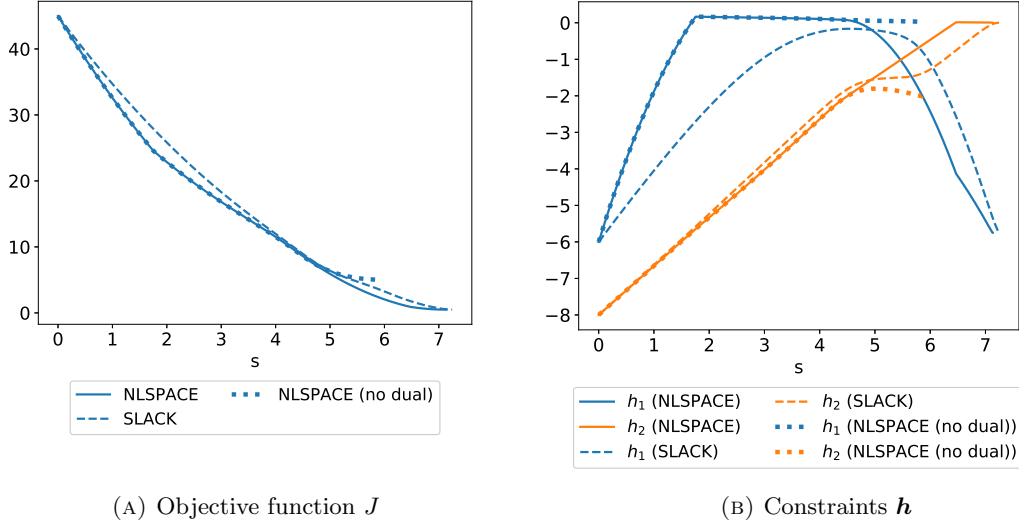


FIGURE 8. History curves of the null space algorithm for the optimization problem of Section 5.3.

- (1) From  $s = 0$  to  $s = 1.73$ , the trajectory  $x(s)$  remains in the feasible domain. Lagrange multipliers  $\mu_1(x(s)) = \mu_2(x(s)) = 0$  are set to 0 and the null space direction  $\xi_J(x(s)) = \nabla J(x(s))$  coincides with the gradient of the objective function, until  $x(s)$  hits the parabolic domain, which corresponds to the saturation of the first constraint  $h_1$ .
- (2) From  $s = 1.73$  to  $s = 4.4$ , the resolution of the dual problem yields a non zero multiplier  $\mu_1(x(s)) > 0$ . The optimization trajectory  $x(s)$  remains tangent to the first constraint, until reaching a limit point such that  $\nabla J(x(s)) \cdot \nabla h_1(x(s)) = 0$ . At this moment, it is not necessary to project the gradient tangentially to this constraint any more, and the values of both Lagrange multipliers  $\mu_1(x(s)) = \mu_2(x(s))$  are equal to 0.
- (3) From  $s = 4.4$  to  $s = 6.5$ , both constraints  $h_1$  and  $h_2$  are ignored and the trajectory  $x(s)$  follows the gradient  $-\nabla J(x(s))$ , till the saturation of  $h_2$ .
- (4) From  $s = 6.5$  to  $s = 7.1$ , the second Lagrange multiplier  $\mu_2(x(s)) > 0$  has a positive value;  $x(s)$  evolves then tangentially to this constraint till the optimum is attained.



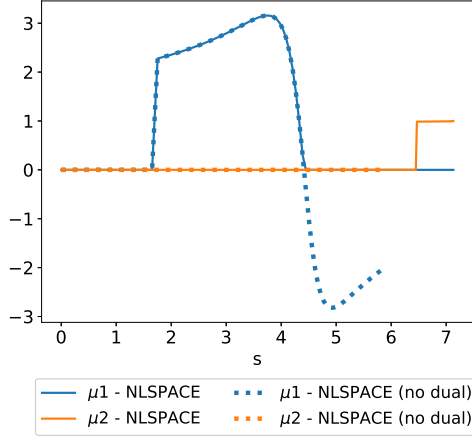


FIGURE 9. Evolution of Lagrange multipliers  $\mu_1(x(s)), \mu_2(x(s))$  for the optimization problem of Section 5.3.

As illustrated on Figure 7, the use of the dual problem is key in the detection of the moment when the optimization trajectory  $x(s)$  needs to be released from active inequality constraints. Because of the discrete nature of the time stepping, the path followed by the ODE (5.1) necessarily enters slightly the violated parabolic domain. Since it does not use the information provided by the dual problem (3.13), the gradient  $\nabla J(x(s))$  is kept being projected tangentially to the constraint  $h_1$  till  $x(s)$  converges to some stationary point (which is not a KKT point). As can be seen on Figure 9, the optimization trajectory followed by this ODE coincides with the one of the flow (3.30), till the instant  $s = 4.4$  at which the Lagrange multiplier  $\mu_1(x(s))$  becomes negative (which violates the feasibility condition of the dual problem (3.13)). Note that using larger steps could have allowed the trajectory  $x(s)$  to exit “by chance” the unfeasible domain, in that case convergence to the optimum would have been obtained, however this would not reflect the actual behavior of the continuous solutions of (5.1).

Finally, the extended ODE (3.8) using slack variables finds a smooth path to the optimum. Since inactive constraints are felt from distance, the trajectory  $x(s)$  is able to remain more strictly in the feasible domain for all times. The total length of the optimization path is almost the same than the one of the null space flow (3.30) (note the steeper descending slopes  $dJ/ds$  for the latter at intersection points of the two trajectories, see the point (3) in the discussion of Section 3.2).

## 6. OPTIMIZATION ON SMOOTH MANIFOLDS: APPLICATION TO SHAPE OPTIMIZATION

As we have already mentioned, our ultimate goal is to apply our optimization strategy to shape and topology optimization problems. Recalling (1.10), the optimization set  $\mathcal{X}$  in this section is therefore a set of shapes in  $\mathbb{R}^d$  ( $d = 2$  or  $3$  in standard applications):

$$\mathcal{X} = \{\Omega \subset D \mid \Omega \text{ Lipschitz}\}, \quad (6.1)$$

where  $D \subset \mathbb{R}^d$  is an enclosing ‘hold-all’ domain. Since  $\mathcal{X}$  is not a Hilbert space, the present context does not fall into the optimization framework described in Sections 2 and 3. However,  $\mathcal{X}$  may be endowed with a manifold structure, which makes it possible to extend our dynamical system (3.9) to this context, up to small adaptations that we now describe.

Building upon the framework of Hadamard’s method of boundary variations, this manifold structure on the set  $\mathcal{X}$  is first defined in Section 6.1. Then, Section 6.2 explains several implementations details of Algorithm 1 that are specific to shape optimization. In particular, we highlight how the classical extension and regularization procedures of shape derivatives are naturally included in our method when using the definition (2.3) of the Hilbertian transposition  $\mathcal{T}$ . Finally, Section 6.3 is devoted to numerical illustrations of our constrained optimization algorithm on the model example of the shape optimization of a bridge structure subjected to multiple load cases.

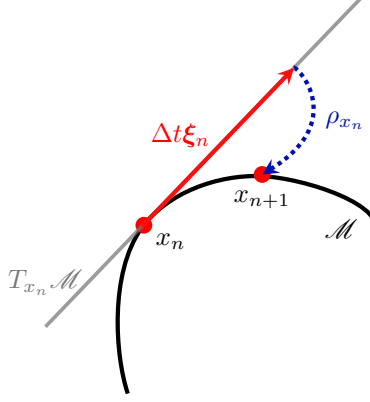


FIGURE 10. Optimization on a manifold  $\mathcal{M}$ : a retraction map  $\rho_{x_n}$  is used to project a tangential motion  $\Delta t \xi_n \in T_{x_n} \mathcal{M}$  from  $x_n \in \mathcal{M}$  back onto the optimization domain  $\mathcal{M}$ .

### 6.1. Hadamard's framework for gradient based shape optimization

Our extension of the previous material to the shape optimization context is inspired by ‘classical’ optimization strategies on a smooth embedded manifold  $\mathcal{M} \subset \mathbb{R}^k$ . In this context, a descent direction at a point  $x_n \in \mathcal{M}$  for some objective functional is typically sought as an element  $\xi_n \in T_{x_n} \mathcal{M}$  of the tangent space  $T_{x_n} \mathcal{M}$  to  $\mathcal{M}$  at  $x_n$ ; see e.g. [24, 1]. Then one relies on a *retraction*  $\rho_{x_n}$ , that is a mapping

$$\rho_{x_n} : T_{x_n} \mathcal{M} \rightarrow \mathcal{M}$$

satisfying the following two consistency conditions:

$$\begin{cases} \rho_{x_n}(0) = x_n \\ \forall \xi \in T_{x_n} \mathcal{M}, \quad \left. \frac{d}{dt} \right|_{t=0} \rho_{x_n}(t\xi) = \xi. \end{cases}$$

The mapping  $\rho_{x_n}$  then allows to convert  $\xi_n$  into a practical update of the actual point  $x_n$  on  $\mathcal{M}$ :

$$x_{n+1} := \rho_{x_n}(\Delta t \xi_n), \quad (6.2)$$

where  $\Delta t > 0$  is the descent step; see [2] and Figure 10. Since the new point  $x_{n+1}$  belongs to  $\mathcal{M}$ , this procedure can be repeated iteratively.

The same idea can be used to apply the methods of Sections 2 and 3 to the optimization problem (1.1), set over the set of shapes  $\mathcal{X}$ . To this end, we rely on Hadamard’s method (see for instance [3, 30, 35, 43]), which considers variations of a shape  $\Omega \in \mathcal{X}$  of the form

$$\rho_\Omega(\theta) := (I + \theta)(\Omega), \text{ for } \theta \in W^{1,\infty}(D, \mathbb{R}^d) \text{ with } \|\theta\|_{W^{1,\infty}(D, \mathbb{R}^d)} < 1. \quad (6.3)$$

Formally, the set  $W^{1,\infty}(D, \mathbb{R}^d)$  may be interpreted as the tangent space of  $\mathcal{X}$  at  $\Omega$  and the mapping  $\rho_\Omega$ , which is defined by (6.3) on a neighborhood of 0 in  $W^{1,\infty}(D, \mathbb{R}^d)$ , plays the role of a retraction. Other definitions are possible for such a transformation dictating how a shape should evolve according to a vector field  $\theta$ , see [4, 20] for discussions regarding that matter. Note also that more rigorous manifold structures on shape spaces can be formulated, see e.g. [9, 39].

Usually, in the context of a general embedded manifold  $\mathcal{M} \subset \mathbb{R}^k$ , a differential structure on  $\mathcal{M}$  is defined first (inducing a notion of derivative on  $\mathcal{M}$ ), and the definition of a suitable retraction is inferred accordingly. In the framework of Hadamard’s method however, it is the retraction  $\rho_\Omega$  itself, that is the parametrization (6.3) by deformation fields  $\theta$ , that is used to define the notion of derivative.

**Definition 5.** A function  $\mathcal{X} \ni \Omega \mapsto F(\Omega) \in \mathbb{R}$  is *shape differentiable* at  $\Omega \in \mathcal{X}$  if the underlying mapping  $F \circ \rho_\Omega : \theta \mapsto F(\rho_\Omega(\theta))$ , from  $W^{1,\infty}(D, \mathbb{R}^d)$  into  $\mathbb{R}$ , is Fréchet differentiable at  $\theta = 0$ . The corresponding derivative

$$DF(\Omega) := D(F \circ \rho_\Omega) : W^{1,\infty}(D, \mathbb{R}^d) \rightarrow \mathbb{R}$$

is called the *shape derivative* of  $F$  at  $\Omega$  and the following expansion holds in the vicinity of  $\boldsymbol{\theta} = 0$ :

$$F(\rho_\Omega(\boldsymbol{\theta})) = F(\Omega) + DF(\Omega)(\boldsymbol{\theta}) + o(\boldsymbol{\theta}), \text{ where } \frac{|o(\boldsymbol{\theta})|}{\|\boldsymbol{\theta}\|_{W^{1,\infty}(D,\mathbb{R}^d)}} \xrightarrow{\boldsymbol{\theta} \rightarrow 0} 0. \quad (6.4)$$

*Remark 9.* In practice, the set  $\mathcal{X}$  of considered shapes is often a subset  $\mathcal{U}_{\text{ad}} \subset \mathcal{X}$  of admissible shapes (e.g. smooth shapes containing some non optimizable regions). The considered deformations are accordingly restrained to a subset  $\Theta_{\text{ad}} \subset W^{1,\infty}(D, \mathbb{R}^d)$ , so that variations of shapes  $\Omega \in \mathcal{U}_{\text{ad}}$  remain in  $\mathcal{U}_{\text{ad}}$ . To keep notations simple, we ignore these details in the presentation; see nevertheless the examples in [Section 6.3](#).

For the shape optimization problem (1.1), we consider objective and constraint functions  $J : \mathcal{X} \rightarrow \mathbb{R}$ ,  $\mathbf{g} : \mathcal{X} \rightarrow \mathbb{R}^p$  and  $\mathbf{h} : \mathcal{X} \rightarrow \mathbb{R}^q$  which are shape differentiable in the sense of [Definition 5](#). Since  $W^{1,\infty}(D, \mathbb{R}^d)$  is not a Hilbert space, the shape derivative  $DJ(\Omega)$  of  $J : \mathcal{X} \rightarrow \mathbb{R}$  (and those of  $\mathbf{g}$  and  $\mathbf{h}$ ) cannot be readily identified with a gradient vector  $\boldsymbol{\xi} \in W^{1,\infty}(D, \mathbb{R}^d)$ . To circumvent this drawback, we introduce a Hilbert space of vector fields  $V \subset W^{1,\infty}(D, \mathbb{R}^d)$ , with inner product  $a(\cdot, \cdot)$ , where the inclusion is continuous. This ensures that  $DJ(\Omega)$ ,  $D\mathbf{g}(\Omega)$  and  $D\mathbf{h}(\Omega)$  are also continuous linear operators on  $V$ , hence the definitions of the gradient  $\nabla J(\Omega) \in V$  and of the transposed operators  $D\mathbf{g}^T(\Omega) : \mathbb{R}^p \rightarrow V$ ,  $D\mathbf{h}^T(\Omega) : \mathbb{R}^q \rightarrow V$  with respect to the inner product  $a$  make sense; see [Definition 1](#). For instance, the gradient  $\nabla J(\Omega) \in V$  is obtained by solving the identification problem:

$$\forall \boldsymbol{\theta} \in V, \quad a(\nabla J(\Omega), \boldsymbol{\theta}) = DJ(\Omega)(\boldsymbol{\theta}). \quad (6.5)$$

Intuitively, the bilinear form  $a$  can be interpreted as a metric on the ‘manifold of shapes’  $\mathcal{X}$ , see e.g. [\[39, 40\]](#). As for the choice of the Hilbert space  $V \subset W^{1,\infty}(D, \mathbb{R}^d)$  used in the identification (6.5), one can for instance take the Sobolev space  $V = H^m(D, \mathbb{R}^d)$  with  $m > 1 + d/2$ , equipped with its standard inner product (the inclusion  $H^m(D, \mathbb{R}^d) \subset W^{1,\infty}(D, \mathbb{R}^d)$  being a consequence of the Sobolev embedding theorem, see [\[15\]](#)). In this case, the identification problem (6.5) boils down to a linear elliptic problem of order  $2m$ .

Let us recall that, under mild regularity assumptions on the objective function  $J(\Omega)$  and the state, the shape derivative of  $J(\Omega)$  can be written in the form of a boundary integral involving only the normal component of the deformation  $\boldsymbol{\theta}$  (this is the so-called *Hadamard structure theorem* [\[30, 35, 42\]](#)). Namely, there exists  $v_J(\Omega) \in L^1(\partial\Omega)$  such that

$$\forall \boldsymbol{\theta} \in W^{1,\infty}(D, \mathbb{R}^d), \quad DJ(\Omega)\boldsymbol{\theta} = \int_{\partial\Omega} v_J(\Omega) \boldsymbol{\theta} \cdot \mathbf{n} ds. \quad (6.6)$$

A very common strategy in the literature (see for instance [\[7, 10, 17, 27, 19, 33\]](#)) consists in taking simply  $H^1(D, \mathbb{R}^d)$  as for the Hilbert space  $V$ , equipped with the inner product

$$\forall \boldsymbol{\theta}, \boldsymbol{\theta}' \in V, \quad a(\boldsymbol{\theta}, \boldsymbol{\theta}') = \int_D (\gamma^2 \nabla \boldsymbol{\theta} : \nabla \boldsymbol{\theta}' + \boldsymbol{\theta} \cdot \boldsymbol{\theta}') dx, \quad (6.7)$$

where  $\gamma > 0$  is a user-defined parameter which can physically be interpreted as a length-scale for the regularity of deformations  $\boldsymbol{\theta}$  (typically,  $\gamma = 3\text{hmin}$  where  $\text{hmin}$  is the minimum edge length of the mesh discretization). Note that this choice of  $V$  is an abuse of the above framework: it is only formal since  $V$  is not a subspace of  $W^{1,\infty}(D, \mathbb{R}^d)$ . However, under the very mild assumption  $v_J(\Omega) \in L^2(\partial\Omega)$ , (which is for instance satisfied in the situations considered in [Section 6.3](#)), the identification problem (6.5) is still well-posed because (6.6) defines a continuous linear form on  $H^1(D, \mathbb{R}^d)$ . In such a situation, the identification (6.5) to (6.7) is interpreted as an extension and regularization of the normal velocity  $v_J(\Omega)$  to the whole domain  $D$ . This step and its consistency with respect to optimization is very classical in shape optimization, see [\[7, 10, 17, 19, 33\]](#). In particular, variants can be considered for tuning more finely the smoothness of such extensions, or to prescribe non optimizable boundaries by prescribing a zero Dirichlet boundary condition in (6.6). Last, the choice  $V = H^1(D, \mathbb{R}^d)$  is quite convenient because this space is easily discretized with  $\mathbb{P}_1$  finite elements. Since this leads to very good results in practice, we shall rely on this strategy in the following.

In light of the previous discussion, the proposed dynamical system (3.9) for tackling shape optimization problems of the form (1.1) is extended and discretized as follows.

- (1) The null space and range space directions  $\boldsymbol{\xi}_J(\Omega)$  and  $\boldsymbol{\xi}_C(\Omega)$  are computed as elements of  $V = H^1(D, \mathbb{R}^d)$  thanks to the formulas (3.11) and (3.29). This requires the computation of the gradient

$\nabla J(\Omega)$  and of the transposes  $D\mathbf{g}^\mathcal{T}(\Omega)$ ,  $D\mathbf{h}^\mathcal{T}(\Omega)$  via the resolution of identification problems such as (6.5). In particular, steps 1 to 4 of Algorithm 1 including the resolution of the dual problem (3.13) are achieved from the knowledge of the Fréchet derivatives and of their transposes.

- (2) The update (3.10) of the design from one iteration of the process to the next is performed by using the retraction map  $\rho_\Omega$  as in (6.2):

$$\Omega_{n+1} := \rho_{\Omega_n}(-\Delta t(\alpha_J \boldsymbol{\xi}_J(\Omega_n) + \alpha_C \boldsymbol{\xi}_C(\Omega_n))); \quad (6.8)$$

the step 6 of Algorithm 1 is adapted accordingly.

The numerical implementation of the retraction map  $\rho_\Omega$ , that is, the procedure used to effectively deform a shape  $\Omega$  according to a prescribed deformation field  $\boldsymbol{\theta}$ , is presented in the next subsection.

## 6.2. Implementation of the constrained gradient flow for level set based shape optimization

Our level set framework for numerical shape and topology optimization is recalled in Section 6.2.1. Further technical details about the practical implementation of Algorithm 1 are then presented in Section 6.2.2.

### 6.2.1. Numerical shape optimization using the level set method and a mesh evolution strategy

Our numerical representation of shapes and their deformations relies on the level set method, pioneered in [37], then introduced in the shape optimization context in [6, 47]. A given shape  $\Omega$  inside the fixed hold-all domain  $D$  is represented by means of a scalar, level set function  $\phi : D \rightarrow \mathbb{R}$  such that:

$$\begin{cases} \phi(x) < 0 & \text{if } x \in \Omega, \\ \phi(x) = 0 & \text{if } x \in \partial\Omega, \\ \phi(x) > 0 & \text{if } x \in D \setminus \overline{\Omega}. \end{cases}$$

The motion of a domain  $\Omega(t)$  in  $D$  evolving over a period of time  $(0, T)$ , starting from a known shape  $\Omega(0) = \Omega$ , according to a vector velocity field  $\boldsymbol{\theta}(x)$  translates in terms of an associated level set function  $\phi(t, x)$  by the following advection equation:

$$\begin{cases} \frac{\partial \phi}{\partial t}(t, x) + \boldsymbol{\theta}(x) \cdot \nabla \phi(t, x) = 0, & t \in (0, T), \quad x \in d, \\ \phi(0, x) = \phi_0(x), & x \in d, \end{cases} \quad (6.9)$$

where  $\phi_0$  is one level set function for  $\omega$ . This approach, which is very convenient for describing the evolution of domains at the discrete level, corresponds actually to the use of a retraction  $\tilde{\rho}_\Omega$  which is slightly different from that  $\rho_\Omega$  defined in (6.3):

$$\tilde{\rho}_\Omega(\boldsymbol{\theta}) := \{x \in D \mid \phi(1, x) < 0\},$$

where  $\phi(t, x)$  is the solution to (6.9). Observe that the mappings  $\tilde{\rho}_\Omega(\boldsymbol{\theta})$  and  $\rho_\Omega(\boldsymbol{\theta})$  differ only from the second order in  $\boldsymbol{\theta}$ ; therefore, the whole optimization process remains consistent in spite of this practical substitution.

In our applications, this method is carried out by relying on the mesh evolution technique of our previous works [5, 27]. At every iteration  $n$ , the current shape  $\Omega_n$  is explicitly discretized as a submesh of a triangulated mesh  $\mathcal{T}_n$  of  $D$  (see e.g. Figure 14 below). A descent direction  $\boldsymbol{\theta}_n(x)$  is computed by estimating

$$\boldsymbol{\theta}_n := -(\alpha_{J,n} \boldsymbol{\xi}_J(\Omega_n) + \alpha_{C,n} \boldsymbol{\xi}_C(\Omega_n)), \quad (6.10)$$

where  $\alpha_J$  and  $\alpha_C$  of the update (6.8) have been replaced by dynamic coefficients  $\alpha_{J,n}$  and  $\alpha_{C,n}$  (this slight modification of Algorithm 1 is detailed in Section 6.2.2 below). A level set function  $\phi_n$  is computed for the current shape  $\Omega_n$ , which is then updated by solving equation (6.9) on the current mesh  $\mathcal{T}_n$  with  $\boldsymbol{\theta} = \boldsymbol{\theta}_n$ :

$$\Omega_{n+1} = \tilde{\rho}_{\Omega_n}(\Delta t \boldsymbol{\theta}_n).$$

In a last stage,  $\mathcal{T}_n$  is remeshed adaptively into a new mesh  $\mathcal{T}_{n+1}$  featuring a discretization of  $\Omega_{n+1}$  as a submesh.

*Remark 10.* In our method, the deformation  $\boldsymbol{\theta}(x)$  is a vector field, in contrast with more classical level set methods [6, 47] that rather rely on a non linear Hamilton-Jacobi equation different to (6.9) involving only the normal component of  $\boldsymbol{\theta}$ . In such settings and more generally, it is convenient to regularize only the normal component  $\boldsymbol{\theta} \cdot \mathbf{n}$  (a scalar field) of the shape derivative since it reduces the size of the identification

problem (6.6). These operations can be achieved consistently with respect to the whole optimization method up to an update of the parameterization space  $V$  and of the inner product  $a$ , see [26] for more details.

### 6.2.2. Determination of the tolerance bounds (4.7) and settings of the parameters $\alpha_{J,n}$ and $\alpha_{C,n}$

We rely on Algorithm 1 for our implementation of the null space flow (3.9) for numerical shape optimization, with the variant introduced in Section 4.2. A few comments are in order regarding the appropriate scaling of the null and range space steps in relation with the size of the mesh discretization; we define accordingly variable coefficients  $\alpha_{J,n}$  and  $\alpha_{C,n}$  for the descent direction  $\boldsymbol{\theta}_n$  in (6.10).

For stability reasons, all vertices of the current shape  $\Omega_n$  should move by only a few mesh elements when reaching the subsequent shape  $\Omega_{n+1}$ . Hence, the minimum edge length  $\mathbf{hmin}$  of the computational mesh is a natural candidate for the limiting step size value  $\mathbf{h}$  of the discussion in Section 4.2. Since *all* values of the displacement  $\boldsymbol{\theta}_n$  should be controlled by  $\mathbf{h}$ , we measure step sizes with the infinity norm  $\|\boldsymbol{\theta}_n\|_{L^\infty(D)}$  rather than with the Hilbertian norm  $\|\boldsymbol{\theta}_n\|_V = \|\boldsymbol{\theta}_n\|_{H^1(D, \mathbb{R}^d)}$  as in Section 4.2. In order to take these changes into account, the tolerance bounds (4.7) need to be updated with respect to this norm as follows:

$$\epsilon_i := \mathbf{h} \int_{\partial\Omega} |v_{C_i}(\Omega_n)| ds,$$

where  $v_{C_i}(\Omega_n) \in L^1(\partial\Omega)$  is the scalar field featured in the shape derivative of the constraint  $C_i(\Omega_n)$  as in (6.6).

Second, the step size  $\|\boldsymbol{\theta}_n\|_{L^\infty(D, \mathbb{R}^d)}$  is effectively controlled by updating dynamically the parameters  $\alpha_{J,n}$  and  $\alpha_{C,n}$  scaling the null space and range space steps  $\boldsymbol{\xi}_J(\Omega_n)$  and  $\boldsymbol{\xi}_C(\Omega_n)$  in (6.10). Since only the products  $\alpha_{J,n}\Delta t$  and  $\alpha_{C,n}\Delta t$  matter, the step size is kept constant and equal to  $\Delta t = 1$ . Then, given  $A_J$  and  $A_C$  two user-defined parameters, which are expressed in terms of the mesh element size  $\mathbf{hmin}$  for a clearer intuitive meaning,  $\alpha_{J,n}$  and  $\alpha_{C,n}$  are updated at every iteration according to the following rules:

$$\alpha_{J,n} := \begin{cases} \frac{A_J \mathbf{hmin}}{\|\boldsymbol{\xi}_J(\Omega_n)\|_{L^\infty(D, \mathbb{R}^d)}} & \text{if } n < \mathbf{n}_0 \\ \frac{A_J \mathbf{hmin}}{\max(\|\boldsymbol{\xi}_J(\Omega_n)\|_{L^\infty(D, \mathbb{R}^d)}, \|\boldsymbol{\xi}_J(\Omega_{\mathbf{n}_0})\|_{L^\infty(D, \mathbb{R}^d)})} & \text{if } n \geq \mathbf{n}_0 \end{cases} \quad (6.11)$$

$$\alpha_{C,n} := \min \left( 0.9, \frac{A_C \mathbf{hmin}}{\max(1\mathbf{e}-9, \|\boldsymbol{\xi}_C(\Omega_n)\|_{L^\infty(D, \mathbb{R}^d)})} \right). \quad (6.12)$$

These normalizations ensure that the null space and range space steps always remain smaller than  $A_J$  and  $A_C$  times the mesh size:

$$\forall n \geq 0, \|\alpha_{J,n} \boldsymbol{\xi}_J(\Omega_n)\|_{L^\infty(D, \mathbb{R}^d)} \leq A_J \mathbf{hmin} \text{ and } \|\alpha_{C,n} \boldsymbol{\xi}_C(\Omega_n)\|_{L^\infty(D, \mathbb{R}^d)} \leq \min(0.9, A_C \mathbf{hmin}).$$

Note that the range step  $\alpha_{C,n} \boldsymbol{\xi}_C(\Omega_n)$  is also set to remain smaller than the constant 0.9, in view of the stability condition  $0 < \alpha_C \Delta t < 2$  (see Remark 8, the role of the constant  $1\mathbf{e}-9$  is only to avoid division by 0 when no constraint is active). In (6.11) and (6.12), the threshold  $\mathbf{n}_0$  is an iteration number indicating for how many steps the normalization by the norm  $\|\boldsymbol{\xi}_J(\Omega_n)\|_{L^\infty(D, \mathbb{R}^d)}$  should be done. By doing so, the null space direction  $\boldsymbol{\xi}_J(\Omega_n)$  is allowed to converge to 0 as  $n \rightarrow \infty$ .

## 6.3. Illustrations on a multiple load structural shape optimization test case

In this final section, we illustrate the efficiency of our optimization strategy on a practical structural design problem. The physical model and the shape optimization setting are presented in Section 6.3.1, then numerical results are discussed for two possible test cases, featuring multiple objective criteria or multiple constraint functions, in Sections 6.3.2 and 6.3.3.

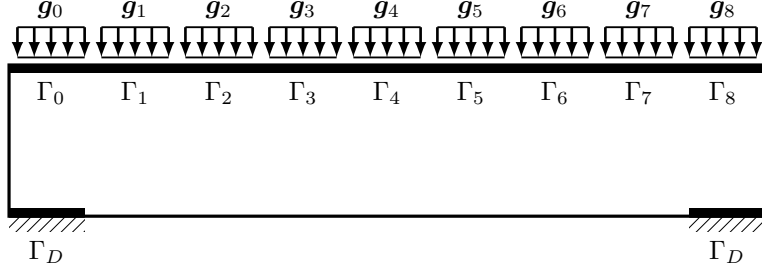


FIGURE 11. Geometric setting for the multiple load case test case

### 6.3.1. Shape optimization setting

We consider the shape optimization of a bridge-like structure  $\Omega$  contained in a two-dimensional rectangular hold-all domain  $D \subset \mathbb{R}^2$  with size  $10 \times 2$ . The boundary of  $\partial\Omega$  is divided into disjoint regions as:

$$\partial\Omega = \Gamma \cup \Gamma_D \cup \bigcup_{i=0}^8 \Gamma_i,$$

where

- $\Gamma_D$  is a non-optimizable part of the boundary on which the structure  $\Omega$  is clamped, made of two segments with unit length at the lower extremities of  $D$ .
- For  $i = 0, \dots, 8$ ,  $\Gamma_i$  is a non-optimizable subset of the upper side of  $D$  with respective abscissa  $[i \frac{10}{9}, (i+1) \frac{10}{9}]$ ;  $\Gamma_i$  is subjected to a unit, vertical downward traction load  $\mathbf{g}_i = (0, -1)$ .
- The remaining region  $\Gamma$  is traction-free and it is the only region of  $\partial\Omega$  which is subject to optimization.

Non-optimizable material layers of width 0.1 are additionally imposed on the upper part of the domain  $D$  and above each component of  $\Gamma_D$ ; see Figure 11. We consider nine different load cases, that are obtained by applying successively and exclusively each of the loads  $\mathbf{g}_i$  on the region  $\Gamma_i$ . In each situation, the corresponding elastic displacement  $\mathbf{u}_i$  is the unique solution in  $H^1(\Omega, \mathbb{R}^d)$  to the linearized elasticity system:

$$\begin{cases} -\operatorname{div}(Ae(\mathbf{u}_i)) = 0 & \text{in } \Omega \\ Ae(\mathbf{u}_i)\mathbf{n} = 0 & \text{on } \Gamma \\ Ae(\mathbf{u}_i)\mathbf{n} = \mathbf{g}_i & \text{on } \Gamma_i \\ Ae(\mathbf{u}_i)\mathbf{n} = 0 & \text{on } \Gamma_j \text{ for } j \neq i \\ \mathbf{u}_i = 0 & \text{on } \Gamma_D, \end{cases} \quad (6.13)$$

where  $e(\mathbf{u}) = (\nabla \mathbf{u} + \nabla \mathbf{u}^T)/2$  is the strain tensor associated to the displacement  $\mathbf{u}$  and  $Ae(\mathbf{u}) = 2\mu e(\mathbf{u}) + \lambda \operatorname{Tr}(e(\mathbf{u}))\mathbf{I}$  is the corresponding stress tensor, involving the Hooke's law  $A$ . The Young modulus and the Poisson ratio are set to  $E = 15$  and  $\nu = 0.35$ , which corresponds to  $\lambda = 12.96$  and  $\mu = 5.56$ . As we have hinted at in Section 6.2.1, the shape is exactly meshed at each iteration (see Figure 14 below), so that each state equation (6.13) is solved by means of a standard finite element method on the meshed subdomain  $\Omega_n$  (without resorting to ersatz material approaches as in e.g. [6]).

Starting from the initial structure  $\Omega_0$  depicted in Figure 12, we minimize the volume  $\operatorname{Vol}(\Omega)$  of the structure  $\Omega$  and maximize the collection of compliances  $C_i(\Omega)$  (for each load case  $\mathbf{g}_i$ ), as a measure of its global rigidity. These quantities are defined by:

$$\operatorname{Vol}(\Omega) := \int_{\Omega} dx, \quad C_i(\Omega) := \int_{\Omega} Ae(\mathbf{u}_i) : e(\mathbf{u}_i) dx, \quad (6.14)$$

and their shape derivatives read (see e.g. [6, 30]):

$$\operatorname{DVol}(\Omega)(\boldsymbol{\theta}) = \int_{\Gamma} \boldsymbol{\theta} \cdot \mathbf{n} ds, \quad \operatorname{DC}_i(\Omega)(\boldsymbol{\theta}) = - \int_{\Gamma} Ae(\mathbf{u}_i) : e(\mathbf{u}_i) \boldsymbol{\theta} \cdot \mathbf{n} ds. \quad (6.15)$$

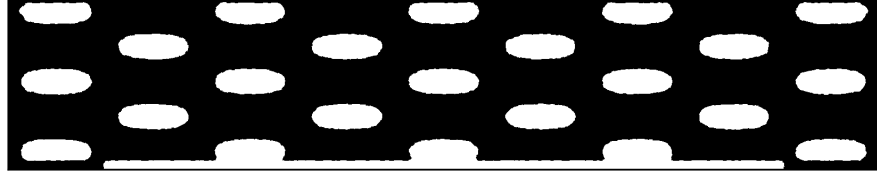


FIGURE 12. Initialisation  $\Omega_0$  (solid in *black*) for the shape optimization examples of [Section 6](#). The thin white layer at the bottom is a non optimizable part of the domain.

### 6.3.2. Volume minimization with maximum compliance constraint

At first, the volume  $\text{Vol}(\Omega)$  is minimized while requiring that each individual compliance  $C_i(\Omega)$  do not exceed a given prescribed value  $C$ :

$$\begin{aligned} \min_{\Omega \in \mathcal{X}} \quad & \text{Vol}(\Omega) \\ \text{s.t.} \quad & C_i(\Omega) \leq C \quad \text{for all } i \in I \end{aligned} \quad (6.16)$$

where  $I \subset \{0, 1, \dots, 8\}$  is a set of indices selecting the considered load cases. We solve (6.16) in the following three configurations:

- (1) *Case 1: single load case:*  $I = \{4\}$  (only the central load  $\mathbf{g}_4$  is applied)
- (2) *Case 2: three load case:*  $I = \{0, 4, 8\}$  (only the central load  $\mathbf{g}_4$  and the two extreme loads  $\mathbf{g}_0$  and  $\mathbf{g}_8$  are applied).
- (3) *Case 3: all load cases:*  $I = \{0, 1, \dots, 8\}$  (all nine loads are considered).

The value of the threshold  $C$  in (6.16) is set to a fraction of the maximum of the compliances  $C_i(\Omega_0)$  of the initial design  $\Omega_0$  (reported on [Figure 12](#)):

$$C = 0.7 \max_{i=0, \dots, 8} \int_{\Omega_0} A e(\mathbf{u}_i) : e(\mathbf{u}_i) dx. \quad (6.17)$$

Let us emphasize that for this example (and the next ones), no fine tuning of the algorithm parameters  $A_J$  and  $A_C$  (determining the update of the values of  $\alpha_{J,n}$  and  $\alpha_{C,n}$  in (6.10)) of [Section 6.2.2](#) is required. The only intuition guiding our choice for this particular test case is that the value of  $A_J$  should be set lower than  $A_C$ . Indeed, a too high value of  $A_J$  might entail a too quick decrease of the volume, which would incur dramatic topological changes violating the rigidity constraints. Therefore these parameters were set to  $A_J = 1$  and  $A_C = 2$  for this test case. The minimum mesh size is  $\text{hmin} = 0.03$ .

The optimized shapes obtained in the three aforementioned situations are shown on [Figure 13](#). The meshes of the initial and final designs, as well as several intermediate shapes corresponding to the nine load test-case are shown on [Figs. 14](#) and [15](#). The convergence histories in the three situations are reported on [Figs. 16](#) to [18](#). They allow to verify the decrease of the objective function even after the saturation of the constraints. Note that for this example and the one to follow, we observed that  $\hat{\Omega}_n$  coincides with  $\tilde{\Omega}_n$  at every iteration, however this situation is very specific to this test case and does not reproduce in generality, see [\[26\]](#) for different shape optimization applications featuring  $\tilde{I}(\Omega_n) \neq \hat{I}(\Omega_n)$ . As expected, the optimum value found for the volume of the solid distribution increases with the number of constraints. The major structural change between the different situations is the addition of extra vertical bars of material near the extremities of the structure.

### 6.3.3. Min/Max compliance optimization with a volume constraint

Now, the maximum value of the compliances  $C_i(\Omega)$  is minimized with an equality volume constraint:

$$\begin{aligned} \min_{\Omega \in \mathcal{X}} \quad & \max_{i \in I} C_i(\Omega) \\ \text{s.t.} \quad & \text{Vol}(\Omega) = \rho_0 \text{Vol}(D) \end{aligned} \quad (6.18)$$



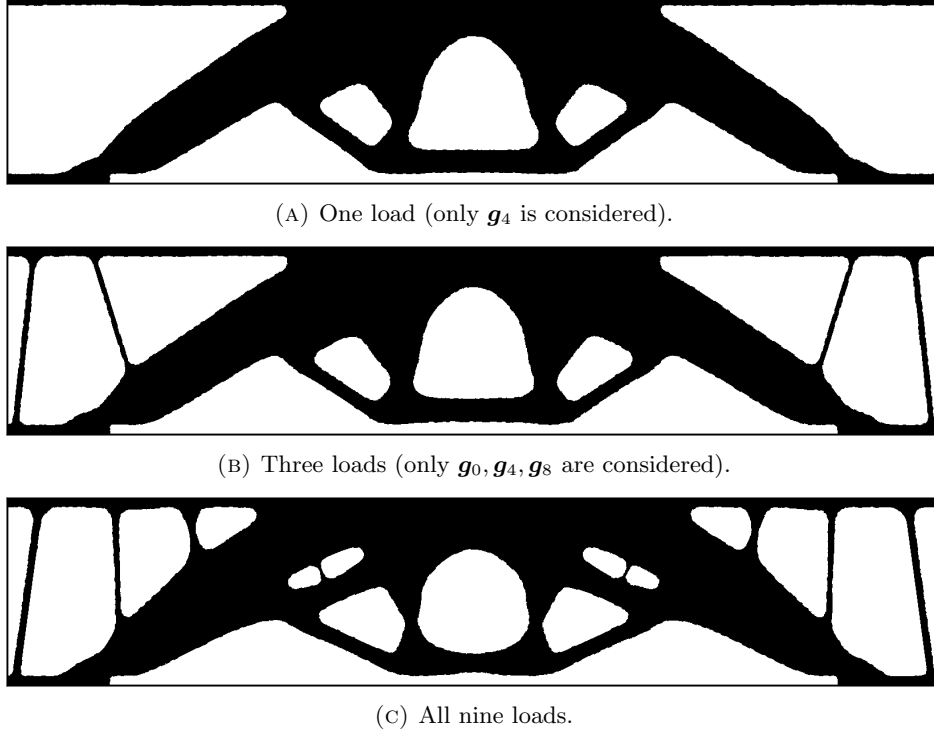


FIGURE 13. Optimized shapes for three possible configurations of the volume minimization problem subject to maximum compliance constraint (Section 6.3.2).

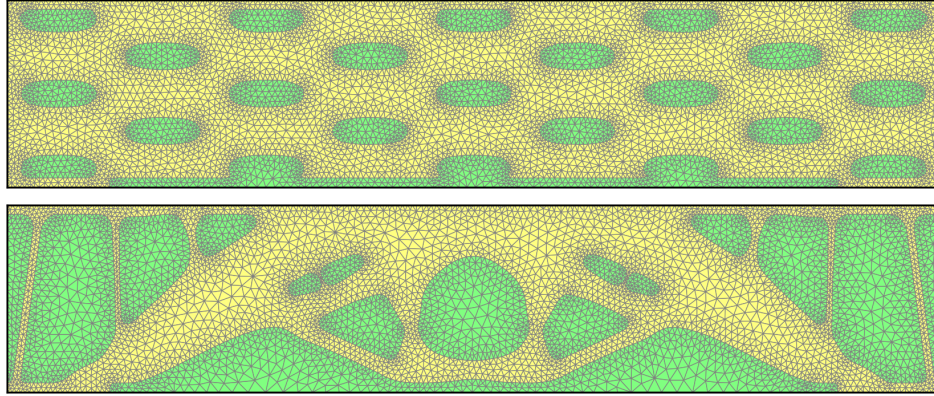


FIGURE 14. Meshes of the initialization and final shapes for the nine load case of Figure 13c (Section 6.3.2).

for a target volume fraction  $\rho_0 = 0.5$  of elastic material and for the three load sets  $I$  introduced in the previous subsection. This problem may be given the form (1.1) after introducing a slack variable  $m$ :

$$\begin{aligned} \min_{(\Omega, m) \in \mathcal{X} \times \mathbb{R}} \quad & m \\ \text{s.t.} \quad & \begin{cases} \text{Vol}(\Omega) = \rho_0 \text{Vol}(D) \\ C_i(\Omega) \leq m \quad \text{for all } i \in I. \end{cases} \end{aligned} \quad (6.19)$$

The optimization must now be performed with respect to both the slack variable  $m$  and the domain geometry  $\Omega$ , which demands minor adaptations of our optimization algorithm (similar e.g. to those in Section 3.2): the

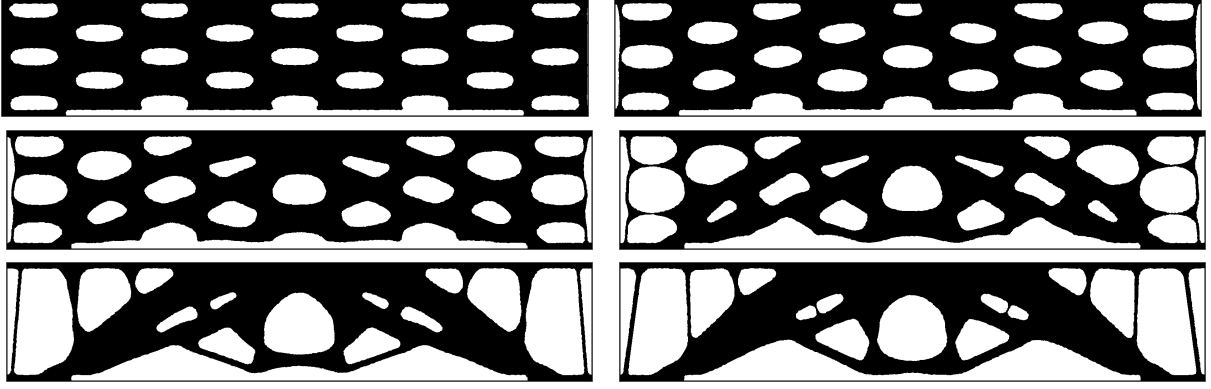


FIGURE 15. Intermediate minimizing shapes for the nine load case of the volume minimization problem of Section 6.3.2 (iterations 0, 5, 10, 20, 80, and 300).

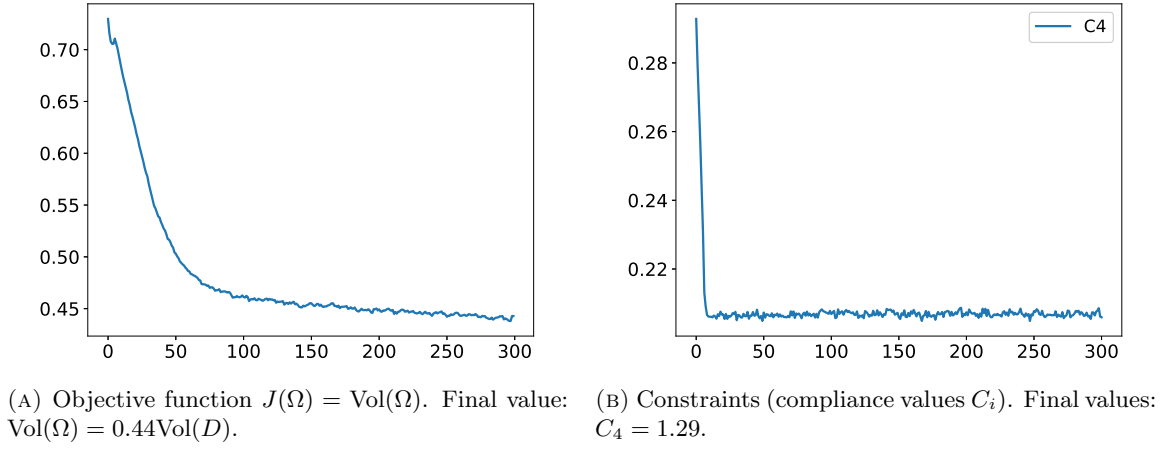


FIGURE 16. Convergence history curves for the single load case of Section 6.3.2.

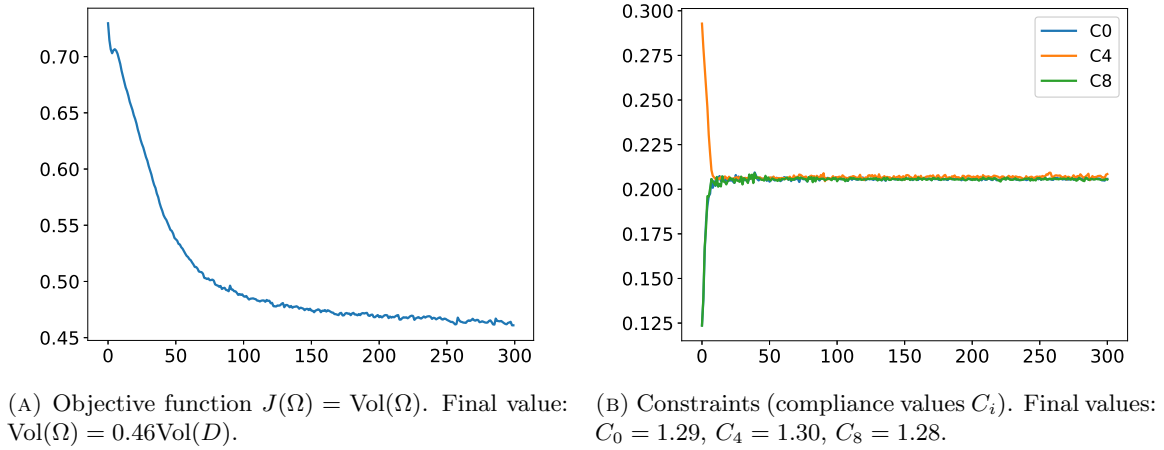


FIGURE 17. Convergence history curves for the three load case of Section 6.3.2.

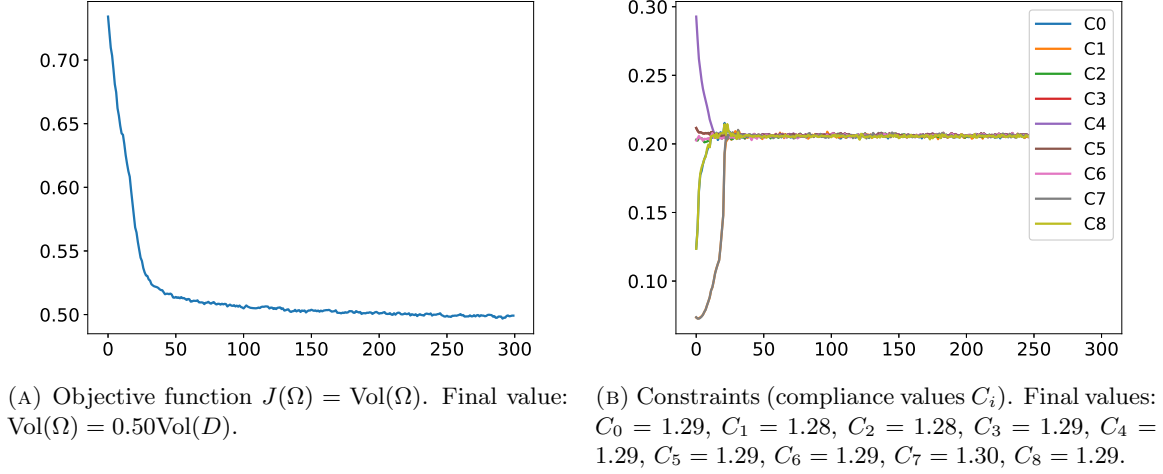


FIGURE 18. Convergence history curves for the nine load case of [Section 6.3.2](#).

optimization domain  $\mathcal{X} \times \mathbb{R}$  is equipped with the tensorized tangent space  $V = H^1(D, \mathbb{R}) \times \mathbb{R}$  and differentials are identified to gradients thanks to the bilinear form  $\tilde{a} : H^1(D, \mathbb{R}) \times \mathbb{R} \rightarrow H^1(D, \mathbb{R}) \times \mathbb{R}$  defined by

$$\forall (v, w) \in H^1(D, \mathbb{R}) \times H^1(D, \mathbb{R}), (l, m) \in \mathbb{R} \times \mathbb{R}, \quad \tilde{a}((v, l), (w, m)) = a(v, w) + lm, \quad (6.20)$$

where  $a$  is the  $H^1(D, \mathbb{R}^d)$  scalar product of (6.7). The slack variable  $m$  is initialized with the maximum value of the compliance of the initial structure  $\Omega_0$  over all the considered loads:

$$m_0 := \max_{i \in I} C_i(\Omega_0), \quad (6.21)$$

and its values  $m_n$  are then updated along with the shape  $\Omega_n$  according to [Algorithm 1](#).

The resulting optimized structures are shown on [Figure 19](#) for each of the three considered configurations and the associated convergence histories are displayed on [Figs. 20 to 22](#) for the single, triple and nine load cases respectively. Note that sudden, abrupt peaks on the constraint curves correspond to topological changes (e.g. at iteration 38 for the nine load case) for which the displacements corresponding to the extremal loads  $\mathbf{g}_0$  and  $\mathbf{g}_8$  are especially sensitive. We observe the decrease of all compliances even after all the inequality constraints are saturated, which occurs as soon as where all compliances achieve a common value. As expected, the optimum design found for the nine load minimum compliance case ([Figure 13c](#)) is similar (up to a few bars) to the corresponding one found for the volume minimization ([Figure 19c](#)): indeed, both cases reach at convergence a volume fraction  $\text{Vol}(\Omega) = 0.5\text{Vol}(D)$  and a maximum compliance  $\max C_i(\Omega) \simeq 1.30$ .

## 7. CONCLUSION AND PERSPECTIVES

In this paper, we have introduced a novel gradient flow for constrained minimization which does not rely on the use of slack variables but on a dual problem allowing to detect on which subset of the constraints the gradient of the objective function should be projected. We insisted on the clear distinction required between differentials and gradients in order to perform such minimization in Hilbert spaces. We provided pedagogical illustrations regarding the behavior of optimization trajectories. Finally, we demonstrated how this flow provides a reliable and comprehensive method for systematic mathematical programming in shape optimization. Further shape optimization applications built upon this algorithm will be provided in the PhD dissertation [\[26\]](#).

**Acknowledgements.** This work was supported by the Association Nationale de la Recherche et de la Technologie (ANRT) [grant number CIFRE 2017/0024]. G. A. is a member of the DEFI project at INRIA Saclay Ile-de-France. The work of C.D. is partially supported by the IRS-CAOS grant from Université Grenoble-Alpes. We thank Alexis Faure for gratefully providing his optimization and plotting routines, which was helpful for making the plots of [Section 5](#).

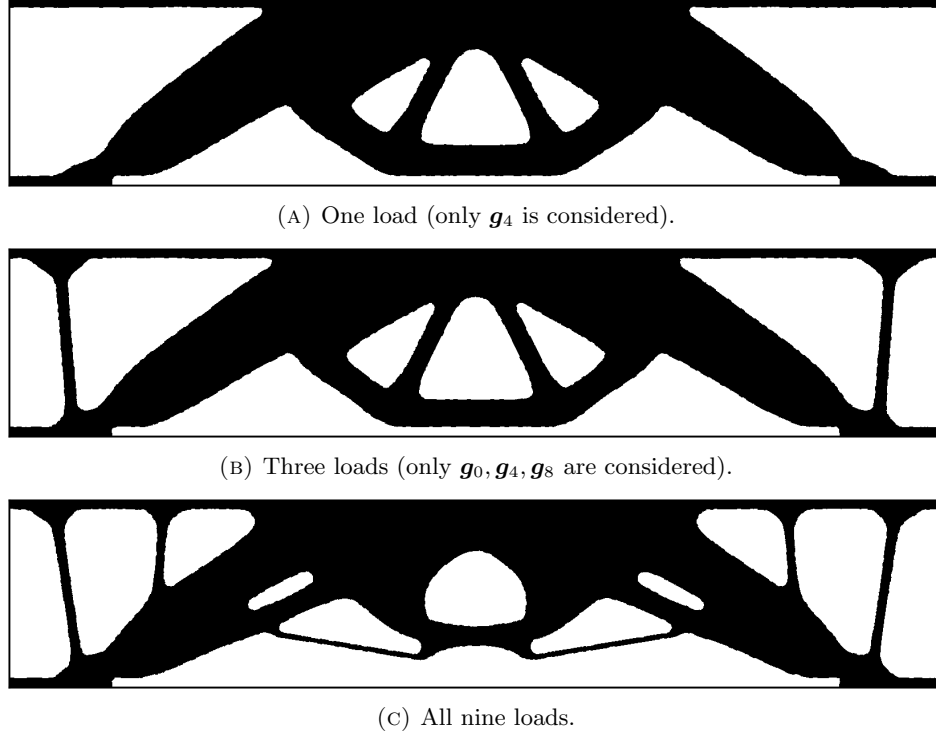


FIGURE 19. Optimized shapes for three possible configurations of the min/max optimization problem (6.19) of Section 6.3.3.

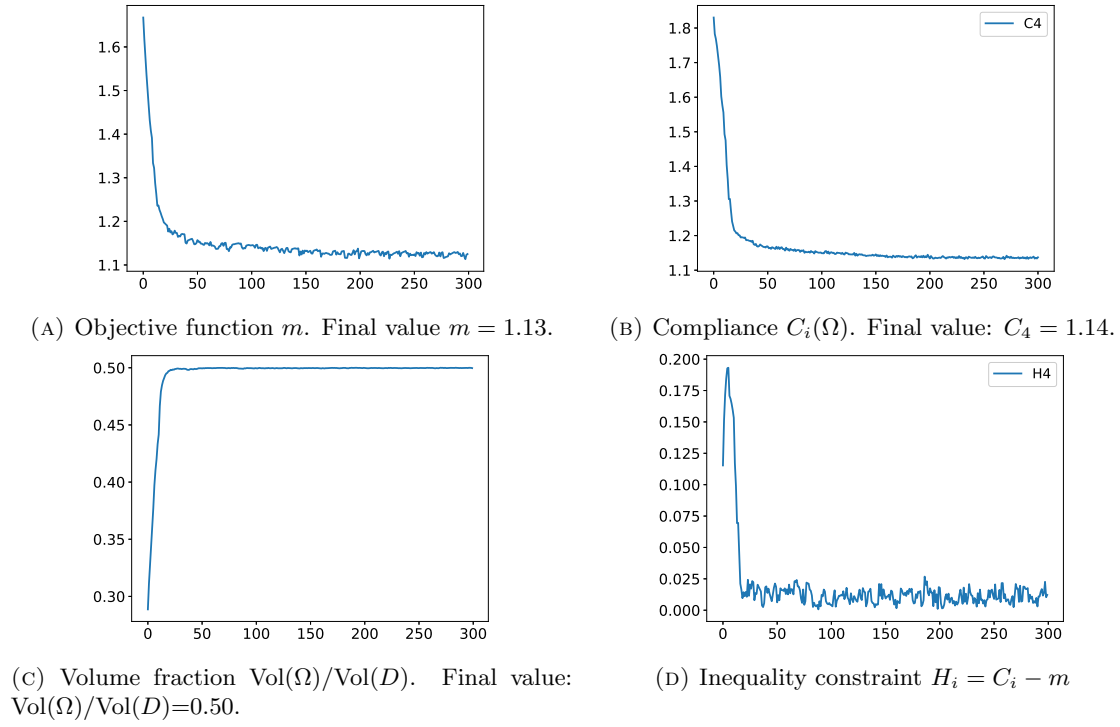


FIGURE 20. Convergence history curves for one load case of Section 6.3.3.

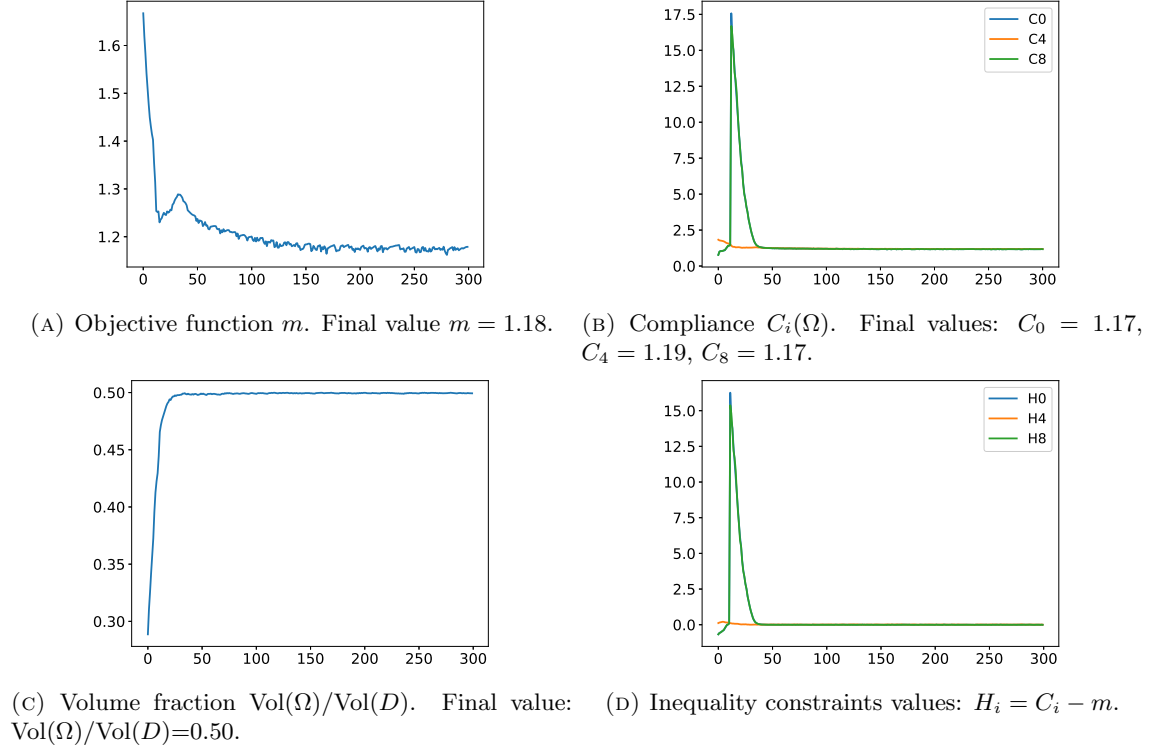


FIGURE 21. Convergence history curves for three load case of Section 6.3.3.

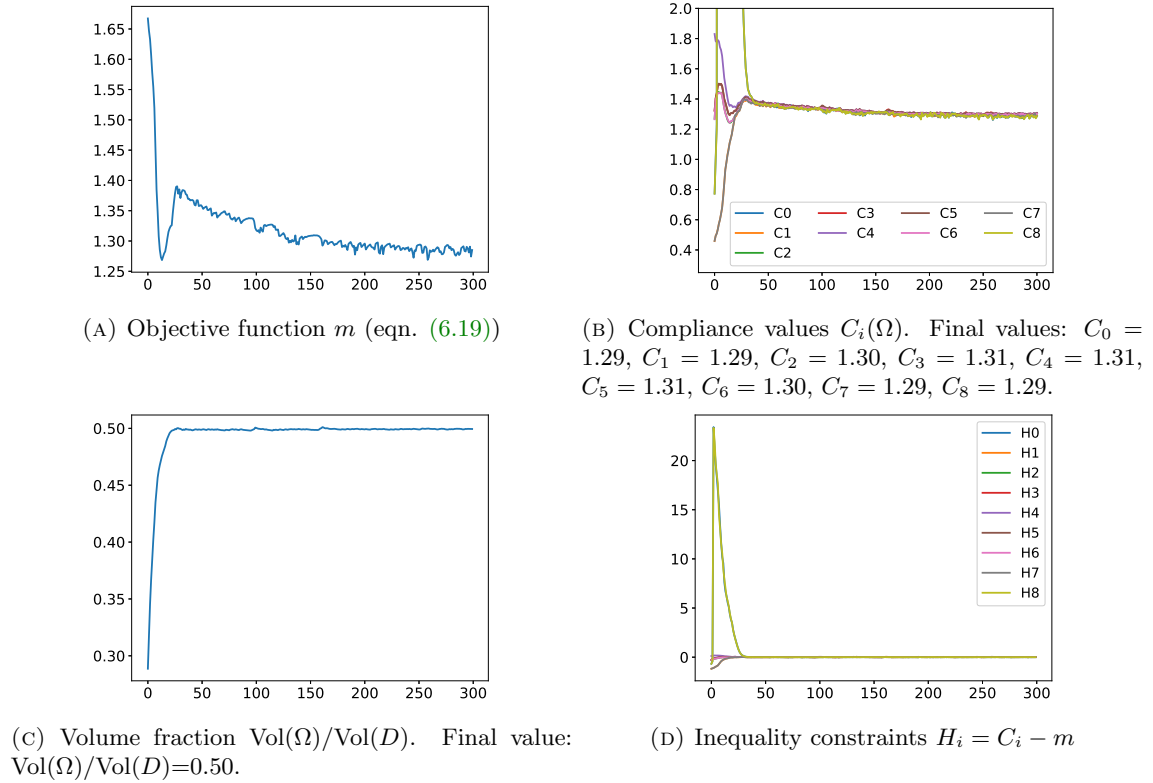


FIGURE 22. Convergence history curves for nine load case of Section 6.3.3.

## REFERENCES

- [1] P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, *Optimization algorithms on matrix manifolds*, Princeton University Press, 2009.
- [2] P.-A. ABSIL AND J. MALICK, *Projection-like retractions on matrix manifolds*, SIAM Journal on Optimization, 22 (2012), pp. 135–158.
- [3] G. ALLAIRE, *Conception optimale de structures, volume 58 of Mathématiques & Applications (Berlin)[Mathematics & Applications]*, Springer-Verlag, Berlin, 2007.
- [4] G. ALLAIRE, E. CANCÈS, AND J.-L. VIÉ, *Second-order shape derivatives along normal trajectories, governed by hamilton-jacobi equations*, Structural and Multidisciplinary Optimization, 54 (2016), pp. 1245–1266.
- [5] G. ALLAIRE, C. DAPOGNY, AND P. FREY, *Shape optimization with a level set based mesh evolution method*, Computer Methods in Applied Mechanics and Engineering, 282 (2014), pp. 22–53.
- [6] G. ALLAIRE, F. JOUVE, AND A.-M. TOADER, *Structural optimization using sensitivity analysis and a level-set method*, Journal of computational physics, 194 (2004), pp. 363–393.
- [7] G. ALLAIRE AND O. PANTZ, *Structural optimization with freefem++*, Structural and Multidisciplinary Optimization, 32 (2006), pp. 173–181.
- [8] M. ANDERSEN, J. DAHL, AND L. VANDENBERGHE, *CVXOPT: A Python package for convex optimization*, Available at <http://cvxopt.org/>, (2012).
- [9] S. ARGUILLÈRE, E. TRÉLAT, A. TROUVÉ, AND L. YOUNES, *Shape deformation analysis from the optimal control viewpoint*, J. Math. Pures Appl. (9), 104 (2015), pp. 139–178.
- [10] H. AZEGAMI AND Z. C. WU, *Domain optimization analysis in linear elastic problems: approach using traction method*, JSME international journal. Ser. A, Mechanics and material engineering, 39 (1996), pp. 272–278.
- [11] C. BARBAROSIE AND S. LOPES, *A gradient-type algorithm for optimization with constraints*, submitted for publication, see also Pre-Print CMAF Pre-2011-001 at <http://cmaf.ptmat.fc.ul.pt/preprints.html>, (2011).
- [12] C. BARBAROSIE, S. LOPES, AND A.-M. TOADER, *A gradient-type algorithm for constrained optimization with applications to multi-objective optimization of auxetic materials*, arXiv preprint arXiv:1711.04863, (2017).
- [13] L. T. BIEGLER, *Nonlinear programming*, vol. 10 of MOS-SIAM Series on Optimization, Society for Industrial and Applied Mathematics (SIAM), 2010.
- [14] J.-F. BONNANS, J. C. GILBERT, C. LEMARÉCHAL, AND C. A. SAGASTIZÁBAL, *Numerical optimization: theoretical and practical aspects*, Springer Science & Business Media, 2006.
- [15] H. BREZIS, *Functional analysis, Sobolev spaces and partial differential equations*, Springer Science & Business Media, 2010.
- [16] R. BRO AND S. DE JONG, *A fast non-negativity-constrained least squares algorithm*, Journal of Chemometrics: A Journal of the Chemometrics Society, 11 (1997), pp. 393–401.
- [17] M. BURGER, *A framework for the construction of level set methods for shape optimization and reconstruction*, Interfaces and Free boundaries, 5 (2003), pp. 301–329.
- [18] C. DAPOGNY, P. FREY, F. OMNÈS, AND Y. PRIVAT, *Geometrical shape optimization in fluid mechanics using FreeFem++*, Structural and Multidisciplinary Optimization, (2017), pp. 1–28.
- [19] F. DE GOURNAY, *Velocity extension for the level-set method and multiple eigenvalues in shape optimization*, SIAM journal on control and optimization, 45 (2006), pp. 343–367.
- [20] M. C. DELFOUR AND J.-P. ZOLÁSIO, *Shapes and geometries: metrics, analysis, differential calculus, and optimization*, vol. 22, Siam, 2011.
- [21] L. DIECI AND L. LOPEZ, *A survey of numerical methods for ivps of odes with discontinuous right-hand side*, Journal of Computational and Applied Mathematics, 236 (2012), pp. 3967–3991.
- [22] J. DIEUDONNÉ, *Foundations of modern analysis*, Academic press, New York and London, 1960.
- [23] P. D. DUNNING AND H. A. KIM, *Introducing the sequential linear programming level-set method for topology optimization*, Structural and Multidisciplinary Optimization, 51 (2015), pp. 631–643.
- [24] A. EDELMAN, T. A. ARIAS, AND S. T. SMITH, *The geometry of algorithms with orthogonality constraints*, SIAM journal on Matrix Analysis and Applications, 20 (1998), pp. 303–353.
- [25] A. FAURE, *Optimisation de forme de matériaux et structures architecturés par la méthode des lignes de niveaux avec prise en compte des interfaces graduées*, PhD thesis, Grenoble Alpes, 2017.
- [26] F. FEPPON, *Numerical methods for shape optimization of multiphysics systems*, PhD thesis, École polytechnique, In preparation.
- [27] F. FEPPON, G. ALLAIRE, F. BORDEU, J. CORTIAL, AND C. DAPOGNY, *Shape Optimization of a Coupled Thermal Fluid-Structure Problem in a Level Set Mesh Evolution Framework*, HAL preprint hal-01686770, (2018).
- [28] A. F. FILIPPOV, *Differential equations with discontinuous righthand sides: control systems*, vol. 18, Springer Science & Business Media, 2013.
- [29] R. FLETCHER, *Practical methods of optimization*, John Wiley & Sons, 2013.
- [30] A. HENROT AND M. PIERRE, *Shape variation and optimization*, vol. 28 of EMS Tracts in Mathematics, European Mathematical Society (EMS), Zürich, 2018. A geometrical analysis, English version of the French publication [MR2512810] with additions and updates.
- [31] M. HINTERMÜLLER, K. ITO, AND K. KUNISCH, *The primal-dual active set strategy as a semismooth newton method*, SIAM Journal on Optimization, 13 (2002), pp. 865–888.
- [32] H. T. JONGEN AND O. STEIN, *Constrained global optimization: adaptive gradient flows*, in Frontiers in global optimization, vol. 74 of Nonconvex Optim. Appl., Kluwer Acad. Publ., Boston, MA, 2004, pp. 223–236.

- [33] B. MOHAMMADI AND O. PIRONNEAU, *Applied shape optimization for fluids*, Oxford University Press, 2010.
- [34] P. MORIN, R. NOCHETTO, M. PAULETTI, AND M. VERANI, *Adaptive sqp method for shape optimization*, in Numerical Mathematics and Advanced Applications 2009, Springer, 2010, pp. 663–673.
- [35] F. MURAT AND J. SIMON, *Sur le contrôle par un domaine géométrique, publications du Laboratoire d'Analyse Numérique*, Université Pierre et Marie Curie, (1976).
- [36] J. NOCEDAL AND S. J. WRIGHT, *Numerical optimization*, Springer Science, 35 (1999).
- [37] S. OSHER AND J. A. SETHIAN, *Fronts propagating with curvature-dependent speed: algorithms based on hamilton-jacobi formulations*, Journal of computational physics, 79 (1988), pp. 12–49.
- [38] J. SCHROPP AND I. SINGER, *A dynamical systems approach to constrained minimization*, Numerical functional analysis and optimization, 21 (2000), pp. 537–551.
- [39] V. H. SCHULZ, *A Riemannian view on shape optimization*, Found. Comput. Math., 14 (2014), pp. 483–501.
- [40] V. H. SCHULZ, M. SIEBENBORN, AND K. WELKER, *Efficient pde constrained shape optimization based on steklov–poincaré-type metrics*, SIAM Journal on Optimization, 26 (2016), pp. 2800–2819.
- [41] V. SHIKHMAN AND O. STEIN, *Constrained optimization: projected gradient flows*, Journal of optimization theory and applications, 140 (2009), pp. 117–130.
- [42] J. SOKOLOWSKI AND J.-P. ZOLESIO, *Introduction to shape optimization*, in Introduction to Shape Optimization, Springer, 1992, pp. 5–12.
- [43] J. SOKOLOWSKI AND J.-P. ZOLESIO, *Introduction to shape optimization*, vol. 16 of Springer Series in Computational Mathematics, Springer-Verlag, Berlin, 1992. Shape sensitivity analysis.
- [44] K. SVANBERG, *The method of moving asymptotes—a new method for structural optimization*, Internat. J. Numer. Methods Engrg., 24 (1987), pp. 359–373.
- [45] K. TANABE, *A geometric method in nonlinear programming*, Journal of Optimization Theory and Applications, 30 (1980), pp. 181–210.
- [46] A. WÄCHTER AND L. T. BIEGLER, *On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming*, Mathematical programming, 106 (2006), pp. 25–57.
- [47] M. Y. WANG, X. WANG, AND D. GUO, *A level set method for structural topology optimization*, Computer methods in applied mechanics and engineering, 192 (2003), pp. 227–246.
- [48] H. YAMASHITA, *A differential equation approach to nonlinear programming*, Mathematical Programming, 18 (1980), pp. 155–168.
- [49] Y.-X. YUAN, *A review of trust region algorithms for optimization*, in ICIAM, vol. 99, Citeseer, 2000, pp. 271–282.
- [50] M. YULIN AND W. XIAOMING, *A level set method for structural topology optimization with multi-constraints and multi-materials*, Acta Mechanica Sinica, 20 (2004), pp. 507–518.
- [51] G. ZOUTENDIJK, *Methods of feasible directions: A study in linear and non-linear programming*, Elsevier Publishing Co., Amsterdam-London-New York-Princeton, N.J., 1960.