

# Discretization-Optimization Methods for Nonlinear Parabolic Relaxed Optimal Control Problems with State Constraints

I. Chrysoverghi<sup>1</sup>, I. Coletsos<sup>1</sup>, J. Geiser<sup>2</sup>, B. Kokkinis<sup>1</sup>

<sup>(1)</sup> Department of Mathematics, School of Applied Mathematics and Physics  
National Technical University of Athens (NTUA)  
Zografou Campus, 15780 Athens, Greece

e-mail: [ichris@central.ntua.gr](mailto:ichris@central.ntua.gr)

<sup>(2)</sup> Weierstrass Institute for Applied Analysis and Stochastics (WIAS)  
Mohrenstrasse 39, D-10117 Berlin, Germany

e-mail: [geiser@wias-berlin.de](mailto:geiser@wias-berlin.de)

## Abstract

We consider an optimal control problem described by a semilinear parabolic partial differential equation, with control and state constraints, where the state constraints and cost involve also the state gradient. Since this problem may have no classical solutions, it is reformulated in the relaxed form. The relaxed control problem is discretized by using a finite element method in space involving numerical integration and an implicit theta-scheme in time for space approximation, while the controls are approximated by blockwise constant relaxed controls. Under appropriate assumptions, we prove that relaxed accumulation points of sequences of optimal (resp. admissible and extremal) discrete relaxed controls are optimal (resp. admissible and extremal) for the continuous relaxed problem. We then apply a penalized conditional descent method to each discrete problem, and also a progressively refining version of this method to the continuous relaxed problem. We prove that accumulation points of sequences generated by the first method are extremal for the discrete problem, and that relaxed accumulation points of sequences of discrete controls generated by the second method are admissible and extremal for the continuous relaxed problem. Finally, numerical examples are given.

**Keywords.** Optimal control, semilinear parabolic systems, state constraints, relaxed controls, discretization,  $\theta$ -scheme, discrete penalized conditional descent method.

**AMS Subject Classification:** 49M25, 49M05, 65N30.

## 1 Introduction

We consider an optimal distributed control problem for systems governed by a semilinear parabolic partial differential equation, with control and state constraints, where the state constraints and cost involve also the gradient of the state. The problem is motivated, for example, by the control of a heat (or other) diffusion process whose source is nonlinear in the heat and temperature, with nonconvex cost and control constraint set (e.g. on-off type control). Since this problem may have no classical solutions, it is reformulated in the relaxed form, using Young measures. The relaxed problem is discretized by using a Galerkin finite element method with continuous piecewise linear basis functions in space and an implicit theta-scheme in time for space approximation, while the controls are approximated by blockwise constant Young measures. We first state the necessary conditions for optimality for the continuous problems, and then for the discrete relaxed problem. Under appropriate assumptions, we prove that relaxed accumulation points of sequences of optimal (resp. admissible and extremal) discrete relaxed controls are optimal (resp. admissible and extremal) for the continuous relaxed problem. We then apply a penalized conditional descent method to each discrete problem, which generates Gamkrelidze

controls, and also a corresponding discretization-optimization method to the continuous relaxed problem that progressively refines the discretization during the iterations, thus reducing computing time and memory. We prove that accumulation points of sequences generated by the fixed discretization method are extremal for the discrete problem, and that relaxed accumulation points of sequences of discrete controls generated by the progressively refining method are admissible and extremal for the continuous relaxed problem. Using a standard procedure, the computed Gamkrelidze controls can then be approximated by piecewise constant classical ones. Finally, numerical examples are given. For approximation of nonconvex optimal control and variational problems, and of Young measures, see e.g. [2-9], [12-14] and the references there.

## 2 The continuous optimal control problems

Let  $\Omega$  be a bounded domain in  $\mathbb{R}^d$ , with boundary  $\Gamma$ , and let  $I = (0, T)$ ,  $T < \infty$ , be an interval. Consider the semilinear parabolic state equation

$$(2.1) \quad y_t + A(t)y + a_0(x, t)^T \nabla y + b(x, t, y(x, t), w(x, t)) = f(x, t, y(x, t), w(x, t))$$

$$(2.2) \quad \text{in } Q = \Omega \times I,$$

$$(2.3) \quad y(x, t) = 0 \text{ in } \Sigma = \Gamma \times I, \quad y(x, 0) = y^0(x) \text{ in } \Omega,$$

where  $A(t)$  is the formal second order elliptic differential operator

$$(2.4) \quad A(t)y := - \sum_{j=1}^d \sum_{i=1}^d (\partial / \partial x_i)[a_{ij}(x, t) \partial y / \partial x_j].$$

The constraints on the control are  $w(x, t) \in U$  in  $Q$ , where  $U$  is a compact, not necessarily convex, subset of  $\mathbb{R}^d$ , the state constraints are

$$(2.5) \quad G_m(w) := \int_Q g_m(x, t, y, \nabla y, w) dx dt = 0, \quad m = 1, \dots, p,$$

$$(2.6) \quad G_m(w) := \int_Q g_m(x, t, y, \nabla y, w) dx dt \leq 0, \quad m = p + 1, \dots, q,$$

and the cost functional to be minimized is

$$(2.7) \quad G_0(w) = \int_Q g_0(x, t, y, \nabla y, w) dx dt.$$

Defining the set of *classical controls*

$$(2.8) \quad W := \{w : (x, t) \mapsto w(x, t) \mid w \text{ measurable from } Q \text{ to } U\},$$

the *continuous classical optimal control problem* is to minimize  $G_0(w)$  subject to  $w \in W$  and to the above state constraints.

Next, we define the set of *relaxed controls* (Young measures; for the relevant theory, see [20], [17])

$$(2.9) \quad R := \{r : Q \rightarrow M_1(U) \mid r \text{ weakly measurable}\} \subset L_w^\infty(Q, M(U)) \equiv L^1(Q, C(U))^*,$$

where  $M(U)$  (resp.  $M_1(U)$ ) is the set of Radon (resp. probability) measures on  $U$ . The set  $R$  is endowed with the relative weak star topology of  $L^1(Q, C(U))^*$ . The set  $R$  is convex, metrizable and compact. If we identify every classical control  $w(\cdot)$  with its associated Dirac relaxed control  $r(\cdot) = \delta_{w(\cdot)}$ , then  $W$  may be considered as a subset of  $R$ , and  $W$  is *thus* dense in  $R$ . For  $\phi \in L^1(Q, C(U)) = L^1(\bar{Q}, C(U))$  (or

$\phi \in B(\bar{Q}, U; \mathbb{R})$ , where  $B(\bar{Q}, U; \mathbb{R})$  is the set of Caratheodory functions in the sense of Warga [20]) and  $r \in L_w^\infty(Q, M(U))$  (in particular, for  $r \in R$ ), we shall use the notation

$$(2.10) \quad \phi(x, t, r(x, t)) := \int_U \phi(x, t, u) r(x, t)(du),$$

and  $\phi(x, t, r(x, t))$  is thus *linear* (under convex combinations, for  $r \in R$ ) in  $r$ . A sequence  $(r_k)$  converges to  $r \in R$  in  $R$  iff

$$(2.11) \quad \lim_{k \rightarrow \infty} \int_Q \phi(x, t, r_k(x, t)) dx dt = \int_Q \phi(x, t, r(x, t)) dx dt,$$

for every  $\phi \in L^1(Q; C(U))$ , or  $\phi \in B(\bar{Q}, U; \mathbb{R})$ , or  $\phi \in C(\bar{Q} \times U)$ .

We denote by  $|\cdot|$  the Euclidean norm in  $\mathbb{R}^n$ , by  $(\cdot, \cdot)$  and  $\|\cdot\|$  the inner product and norm in  $L^2(\Omega)$ , by  $(\cdot, \cdot)_Q$  and  $\|\cdot\|_Q$  the inner product and norm in  $L^2(Q)$ , by  $(\cdot, \cdot)_1$  and  $\|\cdot\|_1$  the inner product and norm in the Sobolev space  $V := H_0^1(\Omega)$ , and by  $\langle \cdot, \cdot \rangle$  the duality bracket between the dual  $V^* = H^{-1}(\Omega)$  and  $V$ . We also define the usual bilinear form associated with  $A(t)$  and defined on  $V \times V$

$$(2.12) \quad a(t, y, v) := \sum_{j=1}^d \sum_{i=1}^d \int_\Omega a_{ij}(x, t) \frac{\partial y}{\partial x_i} \frac{\partial v}{\partial x_j} dx.$$

The relaxed formulation of the above control problem is the following. The relaxed state equation, interpreted in weak form, is

$$(2.13) \quad \langle y_t, v \rangle + a(t, y, v) + (a_0(t)^T \nabla y, v) + (b(t, y, w), v) = (f(t, y, w), v),$$

$\forall v \in V$ , a.e. in  $I$ ,

$$(2.14) \quad y(t) \in V \quad \text{a.e. in } I, \quad y(0) = y^0,$$

(the derivative  $y_t$  is understood here in the sense of  $V$ -vector valued distributions), the control constraint is  $r \in R$ , and the state constraints and cost functionals are

$$(2.15) \quad G_m(r) := \int_Q g_m(x, t, y, \nabla y, r(x, t)) dx dt, \quad m = 0, \dots, q.$$

The *continuous relaxed optimal control problem* is to minimize  $G_0(r)$  subject to the constraints

$$(2.16) \quad r \in R, \quad G_m(r) \leq 0, \quad m = 1, \dots, p, \quad G_m(r) = 0, \quad m = p+1, \dots, q.$$

In the sequel, we shall make some of the following assumptions.

**Assumptions 2.1**  $\Gamma$  is Lipschitz if  $b = 0$ ; else,  $\Gamma$  is  $C^1$  and  $n \leq 3$ .

**Assumptions 2.2** The coefficients  $a_{ij}$  satisfy the ellipticity condition

$$(2.17) \quad \sum_{j=1}^d \sum_{i=1}^d a_{ij}(x, t) z_i z_j \geq \alpha_0 \sum_{i=1}^d z_i^2, \quad \forall z_i, z_j \in \mathbb{R}, \quad (x, t) \in Q,$$

with  $\alpha_0 > 0$ ,  $a_{ij} \in L^\infty(Q)$ , which implies that

$$(2.18) \quad |a(t, y, v)| \leq \alpha_1 \|y\|_1 \|v\|_1, \quad a(t, v, v) \geq \alpha_2 \|v\|_1^2, \quad \forall y, v \in V, \quad t \in I,$$

for some  $\alpha_1 \geq 0$ ,  $\alpha_2 > 0$ .

**Assumptions 2.3**  $a_0 \in L^\infty(Q)^d$ , and the functions  $b, f$  are defined on  $Q \times \mathbb{R} \times U$ , measurable for fixed  $y, u$ , continuous for fixed  $x, t$ , and satisfy the conditions

$$(2.19) \quad |b(x, t, y, u)| \leq \phi(x, t) + \beta |y|^2, \quad b(x, t, y, u) y \geq 0,$$

$$(2.20) \quad |f(x, t, y, u)| \leq \psi(x, t) + \gamma |y|,$$

$$\forall(x, t, y, u) \in Q \times \mathbb{R} \times U,$$

$$(2.21) \quad |f(x, t, y_1, u) - f(x, t, y_2, u)| \leq L|y_1 - y_2|, \quad \forall(x, t, y_1, y_2, u) \in Q \times \mathbb{R}^2 \times U,$$

$$(2.22) \quad b(x, t, y_1, u) \leq b(x, t, y_2, u), \quad \forall(x, t, y_1, y_2, u) \in Q \times \mathbb{R}^2 \times U, \text{ with } y_1 \leq y_2,$$

where  $\phi, \psi \in L^2(Q)$ ,  $\beta, \gamma \geq 0$ .

**Assumptions 2.4** The functions  $g_m$  are defined on  $Q \times \mathbb{R}^{d+1} \times U$ , measurable for fixed  $y, u$ , continuous for fixed  $x, t$ , and satisfy

$$(2.23) \quad |g_m(x, t, y, \bar{y}, u)| \leq \zeta_m(x, t) + \delta_m y^2 + \bar{\delta}_m |\bar{y}|^2, \quad \forall(x, t, y, \bar{y}, u) \in Q \times \mathbb{R}^{d+1} \times U,$$

with  $\zeta_m \in L^1(Q)$ ,  $\delta_m \geq 0$ ,  $\bar{\delta}_m \geq 0$ .

**Assumptions 2.5** The functions  $b, b_y, f_y$  (resp.  $g_m, g_{m_y}, g_{m\bar{y}}$ ) are defined on  $Q \times \mathbb{R} \times U$  (resp.  $Q \times \mathbb{R}^{d+1} \times U$ ), measurable on  $Q$  for fixed  $(y, u) \in \mathbb{R} \times U$  (resp.  $(y, \bar{y}, u) \in \mathbb{R}^{d+1} \times U$ ) and continuous on  $\mathbb{R} \times U$  (resp.  $\mathbb{R}^{d+1} \times U$ ) for fixed  $(x, t) \in Q$ , and satisfy

$$(2.24) \quad |b_y(x, t, y, u)| \leq \xi(x, t) + \eta|y|, \quad |f_y(x, t, y, u)| \leq L_1,$$

$$\forall(x, t, y, u) \in Q \times \mathbb{R} \times U,$$

$$(2.25) \quad |g_{m_y}(x, t, y, \bar{y}, u)| \leq \zeta_{m1}(x, t) + \delta_{m1}|y| + \bar{\delta}_{m1}|\bar{y}|,$$

$$(2.26) \quad |g_{m\bar{y}}(x, t, y, \bar{y}, u)| \leq \zeta_{m2}(x, t) + \delta_{m2}|y| + \bar{\delta}_{m2}|\bar{y}|,$$

$$\forall(x, t, y, \bar{y}, u) \in Q \times \mathbb{R}^{d+1} \times U,$$

with  $\xi, \zeta_{m1}, \zeta_{m2} \in L^2(Q)$ ,  $\eta, \delta_{m1}, \bar{\delta}_{m1}, \delta_{m2}, \bar{\delta}_{m2} \geq 0$ .

The following theorem can be proved by monotonicity and compactness arguments (for continuous  $b, f$  and  $y^0 \in V$ , see also a proof contained in Theorem 3.1 and Lemma 4.2 below).

**Theorem 2.1** Under Assumptions 2.1-3, for every control  $r \in R$  and  $y^0 \in L^2(\Omega)$  (or  $y^0 \in V$ ), the relaxed state equation has a unique solution  $y := y_r$  such that  $y \in L^2(I, V)$ ,  $y_t \in L^2(I, V^*)$ . Moreover,  $y$  is essentially equal to a function in  $C(\bar{I}, L^2(\Omega))$ , and thus the initial condition is well defined.

The following lemma and theorem can be proved by using the techniques of [5], [7], [17].

**Lemma 2.1** Under Assumptions 2.1-3, the operator  $r \mapsto y_r$ , from  $R$  to  $L^2(I, V)$ , and to  $L^2(I, L^4(\Omega))$  if  $b \neq 0$ . Under Assumptions 2.1-4, the functionals  $r \mapsto G_m(r)$ ,  $m = 0, \dots, q$ , from  $R$  to  $\mathbb{R}$ , are continuous.

It is well known that, even if the control set  $U$  is convex, the classical problem may have no classical solutions. But we have anyway the following theorem stating the existence of an optimal relaxed control.

**Theorem 2.2** Under Assumptions 2.1-4, if there exists an admissible control (i.e. satisfying all the constraints), then there exists an optimal relaxed control.

Since  $W \subset R$ , we generally have

$$(2.27) \quad c_R := \min_{\text{constraints on } r} G_0(r) \leq \inf_{\text{constraints on } w} G_0(w) := c_W,$$

where the equality holds, in particular, if there are no state constraints, as  $W$  is dense in  $R$ . Since usually approximation methods slightly violate the state constraints, approximating an optimal relaxed control by a relaxed or a classical control, hence the possibly lower relaxed optimal cost  $c_R$ , is not a drawback in practice (see [20], p. 259).

The following lemma and theorem can be proved by using the techniques of [5], [7], [20] (see also [11]).

**Lemma 2.2** Under Assumptions 2.1-5, dropping the index  $m$  in the functionals, the directional derivative of  $G$  is given, for  $r, r' \in R$ , by

$$(2.28) \quad DG(r, r' - r) := \lim_{\varepsilon \rightarrow 0^+} \frac{G(r + \varepsilon(r' - r)) - G(r)}{\varepsilon} \\ = \int_Q H(x, t, y, \nabla y, z, r'(x, t) - r(x, t)) dx dt,$$

where the Hamiltonian  $H$  is defined by

$$(2.29) \quad H(x, t, y, \bar{y}, z, u) := z[f(x, t, y, u) - b(x, t, y, u)] + g(x, t, y, \bar{y}, u),$$

and the adjoint state  $z = z_r$  satisfies the linear adjoint equation

$$(2.30) \quad -\langle z_r, v \rangle + a(t, v, z) + (a_0^T \nabla v, z) + (z b_y(y, r), v) \\ = (z f_y(y, r) + g_y(y, r), v) + (g_{\bar{y}}(y, \nabla y, r), \nabla v), \quad \forall v \in V, \quad \text{a.e. in } I,$$

$$(2.31) \quad z(t) \in V \quad \text{a.e. in } I, \quad z(T) = 0,$$

with  $y = y_r$ . The mappings  $r \mapsto z_r$ , from  $R$  to  $L^2(Q)$ , and  $(r, r') \mapsto DG(r, r' - r)$ , from  $R \times R$  to  $\mathbb{R}$ , are continuous.

Next, we state the relaxed necessary conditions for optimality.

**Theorem 2.3** Under Assumptions 2.1-5, if  $r \in R$  is optimal for *either* the relaxed or the classical optimal control problem, then  $r$  is *extremal*, i.e. there exist multipliers

$$\lambda_m \in \mathbb{R}, \quad m = 0, \dots, q, \quad \text{with } \lambda_0 \geq 0, \quad \lambda_m \geq 0, \quad m = p+1, \dots, q, \quad \sum_{m=0}^q |\lambda_m| = 1, \quad \text{such that}$$

$$(2.32) \quad \sum_{m=0}^q \lambda_m DG_m(r, r' - r) \geq 0, \quad \forall r' \in R,$$

$$(2.33) \quad \lambda_m G_m(r) = 0, \quad m = p+1, \dots, q \quad (\text{transversality conditions}).$$

The global condition (2.32) is equivalent to the *strong relaxed pointwise minimum principle*

$$(2.34) \quad H(x, t, y(x, t), \nabla y(x, t), z(x, t), r(x, t)) = \min_{u \in U} H(x, t, y(x, t), \nabla y(x, t), z(x, t), u), \\ \text{a.e. in } Q,$$

where complete Hamiltonian and adjoint  $H, z$  are defined with  $g := \sum_{m=0}^q \lambda_m g_m$ .

**Remark.** In the absence of equality state constraints, it can be shown that, if the optimal control  $r$  is *regular*, i.e. there exists  $r' \in R$  such that

$$(2.35) \quad G_m(r) + DG_m(r, r' - r) < 0, \quad m = p+1, \dots, q,$$

(Slater condition), then  $\lambda_0 \neq 0$  for any set of multipliers as in Theorem 2.2.

### 3 The discrete optimal control problems

**Assumptions 3.1**  $\Gamma$  is appropriately piecewise  $C^1$  if  $b=0$ ,  $\Gamma$  is  $C^1$  and  $n \leq 3$  if  $b \neq 0$ ,  $a(t, y, v)$  is independent of  $t$  (for simplicity) and symmetric if  $\theta \neq 1$  in the  $\theta$ -scheme below, the functions  $a_0, b, b_y, b_u, f, f_y, f_u, a_0, b, b_y, f, f_y$  are continuous (possibly finitely piecewise in  $t$ ) on the closure of their domains of definition, and  $y^0 \in V$ .

Under Assumptions 3.1, for each integer  $n \geq 0$ , let  $\Omega^n$  be a subdomain of  $\Omega$  with polyhedral boundary  $\Gamma^n$  such that  $\text{dist}(\Gamma^n, \Gamma) = o(h^n)$ ,  $\{E_i^n\}_{i=1}^{M^n}$  an admissible regular quasi-uniform triangulation of  $\bar{\Omega}^n$  into closed  $d$ -simplices (elements), with  $h^n = \max_i[\text{diam}(E_i^n)] \rightarrow 0$  as  $n \rightarrow \infty$ , and  $\{I_j^n\}_{j=1}^{N^n}$ , a subdivision of the interval  $\bar{I}$  into closed intervals  $I_j^n = [t_{j-1}^n, t_j^n]$ , of equal length  $\Delta t^n$ , with  $\Delta t^n \rightarrow 0$  as  $n \rightarrow \infty$ . We define the *blocks*  $Q_{ij}^n = E_i^n \times I_j^n$ . Let  $V^n \subset V$  be the subspace of functions that are continuous on  $\bar{\Omega}$ , are linear (i.e. affine) on each  $E_i^n$ , and vanish on  $\Omega - \Omega^n$ . Let  $u_0$  be any given fixed point in  $U$ . The set of *discrete classical controls*  $W^n \subset W$  is the subset of classical controls that are constant on the interior of each block  $Q_{ij}^n$  and equal to  $u_0$  on  $Q - (\bar{I} \times \bar{\Omega}^n)$ . The set of *discrete relaxed controls*  $R^n \subset R$  is the subset of relaxed controls that are equal to a constant measure in  $M_1(U)$  on the interior of each block  $Q_{ij}^n$  and equal to  $\delta_{u_0}$  on  $Q - (\bar{I} \times \bar{\Omega}^n)$ . The set  $R^n$  is endowed with the relative weak star topology of  $M(U)^{MN}$ . Clearly, we have  $W^n \subset R^n$ . For implementation reasons, one could alternatively use a coarser partition for the discrete controls, that is, use discrete relaxed controls that are constant on hyperblocks  $Q_{i'j'}^m = E_{i'}^m \times I_{j'}^m$ , where the  $E_{i'}^m$  are appropriate unions of some elements  $E_i^n$  and  $I_{j'}^m$  are appropriate unions of some intervals  $I_j^n$ .

For a given discrete control  $r^n \in R^n$ , and  $\theta \in [1/2, 1]$  if  $b=0$ ,  $\theta=1$  if  $b \neq 0$ , the corresponding discrete state  $y^n := (y_0^n, \dots, y_N^n)$  is given by the discrete state equation (implicit  $\theta$ -scheme)

$$(3.1) \quad (1/\Delta t^n)(y_j^n - y_{j-1}^n, v) + a(y_{j\theta}^n, v) + (a_0^T(t_{j\theta}^n) \nabla y_{j\theta}^n, v) + (b(t_{j\theta}^n, y_{j\theta}^n, r_j^n), v) \\ = (f(t_{j\theta}^n, y_{j\theta}^n, r_j^n), v), \quad \text{for every } v \in V^n, \quad j=1, \dots, N,$$

$$(3.2) \quad (y_0^n - y^0, v)_1 = 0, \quad \text{for every } v \in V^n, \quad y_j^n \in V^n, \quad j=1, \dots, N,$$

where we set

$$(3.3) \quad y_{j\theta}^n := (1-\theta)y_{j-1}^n + \theta y_j^n, \quad t_{j\theta}^n := (1-\theta)t_{j-1}^n + \theta t_j^n.$$

**Theorem 3.1** Under Assumptions 2.2-3 and 3.1, if  $\Delta t^n \leq c'$  (resp.  $\Delta t^n \leq c'(h^n)^2$ ), for some  $c'$  sufficiently small, independent of  $n$  and  $r^n$ , if  $b=0$  (resp.  $b \neq 0$ ), then, for

every  $n$  and every control  $r^n$ , the discrete state equation has a unique solution  $y^n$  such that  $\|y_j^n\| \leq c$ ,  $j=0, \dots, N$ , with  $c$  independent of  $n$  and  $r^n$ .

**Proof (sketch).** Suppose first that  $b \neq 0$ ,  $\theta = 1$ . Lemma 4.1 below shows that, if the solution  $y_j^n$  exists for every  $j$ , then  $\|y_j^n\| \leq c$ ,  $j=0, \dots, N$ , with  $c$  independent of  $n$  and  $r^n$ , for  $\Delta t^n$  as in the assumptions. Suppose by induction that  $y_k^n$  exists for  $k \leq j-1$ . Then the solution  $y_j^n$  is a fixed point of the mapping  $z = F_\theta(y)$  (here with  $\theta = 1$ ), where  $z$  is the solution, for  $y$  given, of the equations

$$(3.4) \quad (1/\Delta t^n)(z - y_{j-1}^n, v) + a(\theta z + (1-\theta)y_{j-1}^n, v) + (a_0^T(t_{j\theta}^n)\nabla(\theta y + (1-\theta)y_{j-1}^n), v) \\ + (b(t_{j\theta}^n, \theta y + (1-\theta)y_{j-1}^n, r_j^n), v) = (f(t_{j\theta}^n, \theta y + (1-\theta)y_{j-1}^n, r_j^n), v), \quad \forall v \in V^n,$$

which reduce (choosing a basis in  $V^n$ ) to a regular linear system in  $z$ . We then show (using our assumptions, the continuous injection  $H_0^1 \subset L^4$ , and the inverse inequality, see [10]) that  $\|z\| = \|F_\theta(y)\| \leq 2c$ , if  $\|y\| \leq 2c$ , for  $\Delta t^n$  as above, i.e.  $F_\theta$  maps the closed ball  $B(0, 2c)$  of center 0 and radius  $2c$  in  $V^n$  into itself. Moreover, one can see (using also the mean value theorem for  $b$ ) that  $F_\theta$  is also contractive in this ball, for  $\Delta t^n$  as above. Therefore  $F_\theta$  has a unique fixed point in  $B(0, 2c)$ , which is the solution  $y_j^n$ . If  $b = 0$ ,  $\theta \in [1/2, 1]$ , one can easily see that, by the Lipschitz continuity of  $f$  in  $y$ , the mapping  $F_\theta$  is contractive on the whole space  $V^n$ , for  $\Delta t^n$  sufficiently small; hence  $F_\theta$  has a unique fixed point in  $V^n$ .

The solution  $y_j^n$  can be computed by the predictor-corrector method, using the linearized semi-implicit predictor scheme, i.e. with  $y_j^{n0} := F_\theta(y_{j-1}^n) \in B(0, 2c)$  or  $V^n$ .

The discrete control constraint is  $r^n \in R^n$  and the discrete functionals are

$$(3.5) \quad G_m^n(r^n) := \Delta t^n \sum_{j=0}^{N-1} \int_{\Omega} g_m(t_{j\theta}^n, y_{j\theta}^n, \nabla y_{j\theta}^n, r_j^n) dx, \quad m = 0, \dots, q.$$

The discrete state constraints are *either* of the *two* following ones

$$(3.6) \quad \text{Case (a)} \quad |G_m^n(r^n)| \leq \varepsilon_m^n, \quad m = 1, \dots, p,$$

$$(3.7) \quad \text{Case (b)} \quad G_m^n(r^n) = \varepsilon_m^n, \quad m = 1, \dots, p,$$

and

$$(3.8) \quad G_m^n(r^n) \leq \varepsilon_m^n, \quad \varepsilon_m^n \geq 0, \quad m = p+1, \dots, q,$$

where the feasibility perturbations  $\varepsilon_m^n$  are given numbers converging to zero, to be defined later. The discrete cost functional to be minimized is  $G_0^n(r^n)$ .

**Theorem 3.2** Under Assumptions 2.2-4 and 3.1, the mappings  $r^n \mapsto y_j^n$  and  $r^n \mapsto G_m^n(r^n)$ , defined on  $R^n$ , are continuous. If any of the discrete problems is feasible, then it has a solution.

**Proof.** The continuity of the operators  $r^n \mapsto y_j^n$  is easily proved by induction on  $j$  (or by using the discrete Bellman-Gronwall inequality, see [18]). The continuity of

$r^n \rightarrow G_m^n(r^n)$  follows from the continuity of  $g_m$ . The existence of an optimal control follows then from the compactness of  $R^n$ .

The proofs of the following lemma and theorem parallel the continuous case and are omitted.

**Lemma 3.2** We drop the index  $m$  in  $g_m$  and  $G_m^n$ . Under Assumptions 2.2-5 and 3.1, for  $r^n, r^m \in R^n$ , the directional derivative of the functional  $G^n$  is given by

$$(3.9) \quad DG^n(r^n, r^m - r^n) = \Delta t^n \sum_{j=0}^{N-1} \int_{\Omega} H(t_{j\theta}^n, y_{j\theta}^n, \nabla y_{j\theta}^n, z_{j,1-\theta}^n, r_j^m - r_j^n) dx,$$

where the discrete adjoint  $z^n$  is given by the linear adjoint scheme

$$(3.10) \quad \begin{aligned} & -(1/\Delta t^n)(z_j^n - z_{j-1}^n, v) + a(v, z_{j,1-\theta}^n) + (a_0^T \nabla v, z_{j,1-\theta}^n) + (z_{j,1-\theta}^n b_y(t_{j\theta}^n, y_{j\theta}^n, r_j^n), v) \\ & = (z_{j,1-\theta}^n f_y(t_{j\theta}^n, y_{j\theta}^n, r_j^n) + g_y(t_{j\theta}^n, y_{j\theta}^n, \nabla y_{j\theta}^n, r_j^n), v) + (g_{\bar{y}}(t_{j\theta}^n, y_{j\theta}^n, \nabla y_{j\theta}^n, r_j^n), \nabla v), \\ & \forall v \in V^n, \quad j = N, \dots, 1, \quad z_N^n = 0, \quad z_j^n \in V^n, \end{aligned}$$

which has a unique solution  $z_{j-1}^n$  for each  $j$ , for  $\Delta t^n$  sufficiently small. Moreover, the mappings  $r^n \mapsto z^n$  and  $(r^n, r^m) \mapsto DG^n(r^n, r^m - r^n)$  are continuous.

**Theorem 3.3** Under Assumptions 2.2-5 and 3.1, if  $r^n \in R^n$  is optimal for the discrete problem with state constraints, Case (b), then it is *extremal*, i.e. there exist multipliers  $\lambda_m^n \in \mathbb{R}$ ,  $m = 0, \dots, q$ , with

$$(3.11) \quad \lambda_m^n \geq 0, \quad \lambda_m^n \geq 0, \quad m = p+1, \dots, q, \quad \sum_{m=0}^q |\lambda_m^n| = 1,$$

such that

$$(3.12) \quad \sum_{m=0}^q \lambda_m^n DG_m^n(r^n, r^m - r^n) = \Delta t^n \sum_{j=1}^N \int_{\Omega} H^n(t_{j\theta}^n, y_{j\theta}^n, \nabla y_{j\theta}^n, z_{j,1-\theta}^n, r_j^m - r_j^n) dx \geq 0,$$

$$\forall r^m \in R^n,$$

$$(3.13) \quad \lambda_m^n [G_m^n(r^n) - \varepsilon_m^n] = 0, \quad m = p+1, \dots, q,$$

where  $H^n$  and  $z^n$  are defined with  $g := \sum_{m=0}^q \lambda_m^n g_m$ . The global condition (3.12) is

equivalent to the *strong discrete blockwise minimum principle*

$$(3.14) \quad \int_{\Omega} H^n(t_j^n, y_{j\theta}^n, \nabla y_{j\theta}^n, z_{j,1-\theta}^n, r_{ij}^n) dx = \min_{u \in U} \int_{\Omega} H^n(t_j^n, y_{j\theta}^n, \nabla y_{j\theta}^n, z_{j,1-\theta}^n, u) dx,$$

$$i = 1, \dots, M, \quad j = 1, \dots, N.$$

## 4 Behavior in the limit

The following control approximation result is proved in [4].

**Proposition 4.1** Under Assumptions 3.1 on  $\Gamma$ , for every  $r \in R$ , there exists a sequence  $(w^n \in W^n)$  that converges to  $r$  in  $R$ .



**Lemma 4.1** (Stability) Under Assumptions 2.2-3 and 3.1, if  $\Delta t$  is sufficiently small, for every  $r^n \in R^n$ , we have the following inequalities, where the constants  $c$  are independent of  $n$  and  $r^n$

$$(4.1) \quad \|y_k^n\| \leq c, \quad k = 0, \dots, N,$$

$$(4.2) \quad \sum_{j=1}^N \|y_j^n - y_{j-1}^n\|^2 \leq c,$$

$$(4.3) \quad \Delta t^n \sum_{j=1}^N \|y_{j\theta}^n\|^2 \leq c,$$

$$(4.4) \quad \Delta t^n \sum_{j=0}^N \|y_j^n\|^2 \leq c, \quad (\text{under the condition } \Delta t^n \leq C(h^n)^2, \text{ for some constant } C \text{ independent of } n, \text{ if } \theta = 1/2),$$

$$(4.5) \quad \Delta t^n \sum_{j=1}^N \|y_j - y_{j-1}\|^2 \leq c, \quad (\text{with the condition } \Delta t^n \leq C(h^n)^2).$$

**Proof.** Dropping the index  $n$  for simplicity of notation, setting  $v = 2\theta\Delta t y_j$  in the discrete equation, and using our assumptions on  $a, a_0, b, f$ , we then have (if  $b \neq 0$ , then  $\theta = 1$  and  $b(y_j)y_j \geq 0$ )

$$(4.6) \quad \begin{aligned} & \theta(\|y_j - y_{j-1}\|^2 + \|y_j\|^2 - \|y_{j-1}\|^2) + (b(t_{j\theta}, y_{j\theta}, r_j), y_j) \\ & + \Delta t [a(y_{j\theta}, y_{j\theta}) + \theta^2 a(y_j, y_j) - (1-\theta)^2 a(y_{j-1}, y_{j-1})] \\ & \leq 2\theta\Delta t |(f(t_{j\theta}, y_{j\theta}, r_j), y_j)| + 2\theta\Delta t \|a_0\|_\infty \|\nabla y_{j\theta}\| \|y_j\| \\ & \leq c\Delta t (1 + \|y_{j\theta}\| + \|y_{j\theta}\|_1) \|y_j\| \\ & \leq c\Delta t [1 + \|y_{j\theta}\|^2 + \beta \|y_{j\theta}\|_1^2 + (1 + \frac{1}{\beta}) \|y_j\|^2] \end{aligned}$$

hence, taking  $\beta \leq \frac{\alpha_2}{2c}$ , we get

$$(4.7) \quad \begin{aligned} & \theta(\|y_j - y_{j-1}\|^2 + \|y_j\|^2 - \|y_{j-1}\|^2) \\ & + \Delta t [\frac{\alpha_2}{2} \|y_{j\theta}\|_1^2 + \theta^2 a(y_j, y_j) - (1-\theta)^2 a(y_{j-1}, y_{j-1})] \\ & \leq c\Delta t (1 + \|y_{j\theta}\|^2 + \|y_j\|^2) \leq c\Delta t (1 + \|y_{j-1}\|^2 + \|y_j\|^2) \\ & \leq c\Delta t (1 + \|y_{j-1}\|^2 + \|y_j - y_{j-1}\|^2), \end{aligned}$$

and if in addition  $\Delta t \leq \frac{\theta}{2c}$

$$(4.8) \quad \begin{aligned} & \theta(\frac{1}{2} \|y_j - y_{j-1}\|^2 + \|y_j\|^2 - \|y_{j-1}\|^2) \\ & + \Delta t [\frac{\alpha_2}{2} \|y_{j\theta}\|_1^2 + \theta^2 a(y_j, y_j) - (1-\theta)^2 a(y_{j-1}, y_{j-1})] \leq c\Delta t (1 + \|y_{j-1}\|^2). \end{aligned}$$

By summation over  $j = 1, \dots, k$ , we obtain, for  $\theta > 1/2$

$$(4.9) \quad \theta(\sum_{j=1}^k \frac{1}{2} \|y_j - y_{j-1}\|^2 + \|y_k\|^2) + \frac{\alpha_2}{2} \Delta t \sum_{j=1}^k \|y_{j\theta}\|_1^2 + \alpha_2 \Delta t c' \sum_{j=1}^k \|y_j\|_1^2$$

$$\leq \theta \|y_0\|^2 + \alpha_1 \Delta t (1-\theta)^2 \|y_0\|_1^2 + c \Delta t \sum_{j=1}^k (1 + \|y_{j-1}\|_1^2), \quad \text{with } c' > 0,$$

and for  $\theta = 1/2$

$$(4.10) \quad \frac{1}{2} \left( \sum_{j=1}^k \frac{1}{2} \|y_j - y_{j-1}\|^2 + \|y_k\|^2 \right) + \frac{\alpha_2}{2} \Delta t \sum_{j=1}^k \|y_{j\theta}\|_1^2 \\ \leq \frac{1}{2} \|y_0\|^2 + \alpha_1 \frac{\Delta t}{4} \|y_0\|_1^2 + c \Delta t \sum_{j=1}^k (1 + \|y_{j-1}\|_1^2).$$

Since  $\|y_0\|_1$ , hence  $\|y_0\|$ , remains bounded, using the discrete Bellman-Gronwall inequality (see [18]), we obtain inequality (i). The inequalities (4.2), (4.3), and (4.4) if  $\theta > 1/2$ , follow. By the inverse inequality (see [10]), the condition  $\Delta t^n \leq C(h^n)^2$ , and inequality (4.2), we get inequality (4.5)

$$(4.11) \quad \Delta t \sum_{j=1}^N \|y_j - y_{j-1}\|_1^2 \leq \frac{\Delta t}{h^2} \sum_{j=1}^N \|y_j - y_{j-1}\|^2 \leq C \sum_{j=1}^N \|y_j - y_{j-1}\|^2 \leq c.$$

If  $\theta = 1/2$ , inequality (4.4) follows from inequalities (4.3) and (4.5).

For given values  $v_0, \dots, v_N$  in a vector space, we define the piecewise constant and continuous piecewise linear functions

$$(4.12) \quad v_-(t) := v_{j-1}, \quad v_+(t) := v_j, \quad v_\theta(t) := (1-\theta)v_{j-1} + \theta v_j, \quad t \in \overset{o}{I}_j^n, \quad j = 1, \dots, N,$$

$$(4.13) \quad v_\wedge(t) := v_{j-1} + \frac{t - t_{j-1}^n}{\Delta t^n} (v_j - v_{j-1}), \quad t \in I_j^n, \quad j = 1, \dots, N.$$

If  $b = 0$  (resp.  $b \neq 0$ ), we suppose in the sequel that  $\Delta t^n \leq C$  (resp.  $\Delta t^n \leq C(h^n)^2$ ), with  $C$  sufficiently small, so as to guarantee the results of Theorem 3.1 and Lemma 4.1.

**Lemma 4.2** (Consistency of states and functionals) Under Assumptions 2.2-3 and 3.1, if  $r^n \rightarrow r$  in  $R$ , then the corresponding discrete states  $y_\wedge^n, y_+^n, y_-^n, y_\theta^n$  converge to  $y_r$  in  $L^2(I, L^4(\Omega))$  (resp  $L^2(Q)$ ) strongly if  $b \neq 0$  (resp.  $b = 0$ ),  $y_\theta^n \rightarrow y_r$  in  $L^2(I, V)$  strongly, and

$$(4.14) \quad \lim_{n \rightarrow \infty} G_m^n(r^n) = G_m(r), \quad m = 0, \dots, q.$$

**Proof.** By Lemma 4.1 (inequality (4.2) multiplied by  $\Delta t$ ),  $y_+^n - y_-^n \rightarrow 0$  in  $L^2(Q)$  strongly. Since, by inequality (4.3) in Lemma 3.3,  $y_-^n$  and  $y_+^n$  are bounded in  $L^2(I, V)$ , it follows that  $y_\wedge^n$  is also bounded in  $L^2(I, V)$ . By extracting subsequences, we can suppose that  $y_\wedge^n \rightarrow y$  and  $y_\theta^n \rightarrow y$  in  $L^2(I, V)$  weakly (hence in  $L^2(Q)$  weakly), for the same  $y$ . The discrete state equation can be written in the form

$$(4.15) \quad \frac{d}{dt} (y_\wedge^n(t), v^n) = (\psi^n(t), v^n)_1, \quad \forall v^n \in V^n, \text{ a.e. in } (0, T),$$

in the scalar distribution sense, where the piecewise constant function  $\psi^n$  is defined, using Riesz's representation theorem, by

$$(4.16) \quad (\psi_j^n(t), v^n)_1 := -a(y_{j\theta}^n, v^n) - (a_0(t_{j\theta}^n)^T \nabla y_{j\theta}^n, v^n) - (b(t_{j\theta}^n, y_{j\theta}^n, r_j^n), v^n) \\ + (f(t_{j\theta}^n, y_{j\theta}^n, r_j^n), v^n), \quad \text{in } \overset{o}{I}_j^n, \quad j = 1, \dots, N.$$

By our assumptions, we have, for  $j = 1, \dots, N$

$$(4.17) \quad \begin{aligned} |(\psi_j^n, v^n)_1| &\leq c(\|y_{j\theta}^n\|_1 \|v^n\|_1 + \sigma(1 + \|y_{j\theta}^n\|_{L^4(\Omega)}^2) \|v^n\| + (1 + \|y_{j\theta}^n\|_1) \|v^n\|) \\ &\leq c(1 + \sigma \|y_{j\theta}^n\|_1^2 + \|y_{j\theta}^n\|_1) \|v^n\|_1, \end{aligned}$$

with  $\sigma = 0$  if  $b = 0$ ,  $\sigma = 1$  if  $b \neq 0$ ; hence

$$(4.18) \quad \|\psi_j^n\|_1 \leq c(1 + \|y_{j\theta}^n\|_1^2 + \|y_{j\theta}^n\|_1) \leq c(1 + \|y_{j\theta}^n\|_1^2).$$

Therefore, using inequality (4.3) in Lemma 3.3

$$(4.19) \quad \int_0^T \|\psi^n(t)\|_1 dt \leq c(1 + \int_0^T \|y_\theta^n\|_1^2 dt) \leq c,$$

which shows that  $\psi^n$  belongs to  $L^1(I, V)$ . Now, Let  $\tilde{\psi}^n$  denote the extension of  $\psi^n$  by 0 outside  $[0, T]$ . We then have, on  $\mathbb{R}$

$$(4.20) \quad \frac{d}{dt} (y_\lambda^n(t), v^n) = (\psi^n(t), v^n)_1 + (y_0^n, v^n) \delta_0 - (y_N^n, v^n) \delta_T,$$

where  $\delta_0, \delta_T$  are the Dirac distributions at 0 and  $T$ . Taking the Fourier transforms

( $\hat{\psi}^n$  Fourier transform of  $\tilde{\psi}^n$ ), we have

$$(4.21) \quad -2i\pi\tau (\hat{y}_\lambda^n, v^n) = (\hat{\psi}^n(\tau), v^n) + (y_0^n, v^n) - (y_N^n, v^n) e^{-2i\pi\tau T}.$$

Setting  $v^n = \hat{y}_\lambda^n(\tau)$ , and since  $y_0^n, y_N^n$  are bounded in  $L^2(\Omega)$ , we get

$$(4.22) \quad 2\pi|\tau| \|\hat{y}_\lambda^n(\tau)\|^2 \leq \|\hat{\psi}^n(\tau)\|_1 \|\hat{y}_\lambda^n(\tau)\|_1 + c \|\hat{y}_\lambda^n(\tau)\|.$$

By the definition of the Fourier transform, we obtain

$$(4.23) \quad \|\hat{\psi}^n(\tau)\|_1 \leq \int_0^T \|\psi^n(t)\|_1 dt \leq c.$$

Therefore

$$(4.24) \quad |\tau| \|\hat{y}_\lambda^n(\tau)\|^2 \leq c \|\hat{y}_\lambda^n(\tau)\|_1.$$

For  $\rho \in [0, 1/4)$ , the following inequality holds on  $\mathbb{R}$

$$(4.25) \quad |\tau|^{2\rho} \leq c \frac{1 + |\tau|}{1 + |\tau|^{1-2\rho}}.$$

We then have

$$(4.26) \quad \begin{aligned} \int_{-\infty}^{+\infty} |\tau|^{2\rho} \|\hat{y}_\lambda^n(\tau)\|^2 d\tau &\leq c \int_{-\infty}^{+\infty} \frac{1 + |\tau|}{1 + |\tau|^{1-2\rho}} \|\hat{y}_\lambda^n(\tau)\|^2 d\tau \\ &\leq c \int_{-\infty}^{+\infty} \|\hat{y}_\lambda^n(\tau)\|^2 d\tau + c' \int_{-\infty}^{+\infty} \frac{\|\hat{y}_\lambda^n(\tau)\|_1}{1 + |\tau|^{1-2\rho}} d\tau \\ &\leq c \int_{-\infty}^{+\infty} \|\hat{y}_\lambda^n(\tau)\|^2 d\tau + c' \left[ \int_{-\infty}^{+\infty} \frac{d\tau}{(1 + |\tau|^{1-2\rho})^2} \right]^{1/2} \left( \int_{-\infty}^{+\infty} \|\hat{y}_\lambda^n(\tau)\|_1^2 d\tau \right)^{1/2}. \end{aligned}$$

The constant integral factor here is finite for  $\rho < 1/4$ . By the Parseval identity

$$(4.27) \quad \int_{-\infty}^{+\infty} \|\hat{y}_\lambda^n(\tau)\|^2 d\tau \leq c \int_{-\infty}^{+\infty} \|\hat{y}_\lambda^n(\tau)\|_1^2 d\tau = c \int_0^T \|y_\lambda^n(t)\|_1^2 dt \leq c.$$

Therefore, we obtain

$$(4.28) \quad \int_{-\infty}^{+\infty} |\tau|^{2\rho} \|\hat{y}_\lambda^n(\tau)\|^2 d\tau \leq c.$$

Let us examine first the case  $b \neq 0$  with its assumptions. By the Compactness Theorem 2.2, Ch. III, in [18], and since the injection of  $V = H_0^1(\Omega)$  into  $H^{1-\varepsilon}(\Omega)$ ,

$\varepsilon \in (0, 1]$ , is compact, there exists a subsequence (same notation) such that  $y_\wedge^n \rightarrow \tilde{y}$  in  $L^2(I, H^{1-\varepsilon}(\Omega))$  and in  $L^2(Q)$  strongly, for some  $\tilde{y}$ , and we must have  $\tilde{y} = y$ , since  $\hat{y}_\wedge^n \rightarrow y$  also in  $L^2(Q)$  weakly. Since the injection of  $H^{1-\varepsilon}(\Omega)$  into  $L^4(\Omega)$  is continuous for  $\varepsilon$  sufficiently small (see [1]),  $y_\wedge^n \rightarrow y$  in  $L^2(I, L^4(\Omega))$  strongly. Since, by Lemma 4.1 (v),  $y_+^n - y_-^n \rightarrow 0$  in  $L^2(I, V)$ , hence in  $L^2(I, L^4(\Omega))$ , strongly, we get also that  $y_-^n, y_+^n \rightarrow y$  in  $L^2(I, L^4(\Omega))$  strongly. In the case  $b = 0$ , we can use directly the compactness of the injection  $V \subset L^2(\Omega)$  (between Hilbert spaces) to show the strong convergences  $y_\wedge^n, y_-^n, y_+^n, y_\theta^n \rightarrow y$  in  $L^2(Q)$ , using inequality (4.2) (instead of (4.5)) in Lemma 4.1. Now, to show that  $y = y_r$ , we proceed similarly to the proof of Lemma 4.3 in [5], i.e. we pass to the limit in the discrete equation, integrated in  $t$ , with appropriate interpolating test functions  $\phi^n(t)v^n(x)$  (with  $v^n$  modified near the boundary so as to belong to  $V^n$ ); for the passage to the limit in the nonlinear terms containing  $b$  and  $f$ , we can use a generalization of Proposition 2.1 in [4] for a double integral whose integrand involves multiple sequences converging in various  $L^p$  spaces, which can be proved by using the convergence  $r^n \rightarrow r$  in  $R$ , the fact that converging sequences of functions in  $L^p$  are dominated (in norm a.e. in  $\Omega$ , and up to subsequences) by a fixed function in  $L^p$ , Hölder's inequality, Egorov's theorem, and Lebesgue's dominated convergence theorem. Next, to prove the strong convergence  $y_\theta^n \rightarrow y$  in  $L^2(I, V)$ , we first remark that, by the discrete and continuous state equations, the boundedness of  $(y_N^n)$  in  $L^2(\Omega)$  by inequality (4.1) in Lemma 4.1, the above convergences, Proposition 2.1 in [4], and taking the sequence  $(v^n \in V^n)$  of functions interpolating an arbitrary given  $v \in C_0^1(\bar{\Omega})$  at the vertices inside  $\Omega^n$  and vanishing on  $\Gamma^n$  (which converges to  $v$  in  $V$  strongly), we have

$$\begin{aligned}
(4.29) \quad & (y_N^n, v) = (y_N^n, v - v^n) + (y_N^n, v^n) \\
& = (y_N^n, v - v^n) + (y_0^n, v^n) + \int_0^T (f(t_\theta^n, y_\theta^n, r^n), v^n) dt \\
& \quad - \int_0^T (b(t_\theta^n, y_\theta^n, r^n), v^n) dt - \int_0^T (a_0^T \nabla y_\theta^n, v^n) dt - \int_0^T a(y_\theta^n, v^n) dt \\
& \rightarrow (y^0, v) + \int_0^T (f(y, r), v) dt - \int_0^T (b(y, r), v) dt - \int_0^T (a_0^T \nabla y, v) dt - \int_0^T a(y, v) dt \\
& = (y(T), v),
\end{aligned}$$

for every  $v \in C_0^1(\bar{\Omega})$ , hence  $(y_N^n, v) \rightarrow (y(T), v)$  for every  $v \in L^2(\Omega)$ , since  $C_0^1(\bar{\Omega})$  is dense in  $L^2(\Omega)$ , i.e.  $y_N^n \rightarrow y(T)$  in  $L^2(\Omega)$  weakly. We then have

$$\begin{aligned}
(4.30) \quad & \alpha_2 \|y_+^n - y\|_{L^2(I, V)}^2 \leq \int_0^T a(y_\theta^n - y, y_\theta^n - y) dt + \frac{1}{2} \|y_N^n - y(T)\|^2 \\
& = \frac{1}{2} \|y_0^n\|^2 - \frac{1}{2} (y_N^n, y(T)) - \frac{1}{2} (y(T), y_N^n - y(T)) \\
& \quad + \int_0^T (f(y_\theta^n, r^n), y_\theta^n) dt - \int_0^T (b(y_\theta^n, r^n), y_\theta^n) dt - \int_0^T (a_0^T \nabla y_\theta^n, y_\theta^n) dt \\
& \quad - \int_0^T a(y_\theta^n, y) dt - \int_0^T a(y, y_\theta^n - y) dt \rightarrow 0.
\end{aligned}$$

Finally, the last convergences of the lemma follow from Proposition 2.1 in [4].

Note that the condition  $\Delta t^n \leq C(h^n)^2$  imposed (in fact, the inverse inequality used to derive inequality (v)) is a worst case one. In practice, the corresponding sequences of gradients  $(\nabla y^n)$  constructed by the algorithms are often bounded in  $L^2(Q)$ , or even in  $L^\infty(Q)$ , and the above condition is not needed.

We suppose in the sequel that the continuous relaxed problem is feasible. The following (theoretical, in the presence of state constraints) theorem addresses the behavior in the limit of optimal discrete controls.

**Theorem 4.1** We suppose that Assumptions 2.2-4 and 3.1 are satisfied. In the presence of state constraints, we suppose in addition that the sequences  $(\varepsilon_m^n)$  in the discrete state constraints, Case (a), converge to zero as  $n \rightarrow \infty$  and satisfy

$$|G_m^n(\tilde{r}^n)| \leq \varepsilon_m^n, \quad m = 1, \dots, p, \quad G_m^n(\tilde{r}^n) \leq \varepsilon_m^n, \quad \varepsilon_m^n \geq 0, \quad m = p+1, \dots, q,$$

for every  $n$ , where  $(\tilde{r}^n \in R^n)$  is a sequence converging in  $R$  to an optimal control  $\tilde{r} \in R$  of the relaxed problem. For each  $n$ , let  $r^n$  be optimal for the discrete problem, Case (a). Then every relaxed accumulation point of  $(r^n)$  is optimal for the continuous relaxed problem.

**Proof.** Note that our assumption implies that the discrete problems are feasible for every  $n$ . Let  $(r^n)$  be a subsequence (same notation) that converges to some  $r \in R$ . Since  $r^n$  is optimal, hence admissible, and  $\tilde{r}^n$  is admissible, for the discrete problem, we have

$$(4.31) \quad G_0^n(r^n) \leq G_0^n(\tilde{r}^n), \quad |G_m^n(r^n)| \leq \varepsilon_m^n, \quad m = 1, \dots, p, \quad G_m^n(r^n) \leq \varepsilon_m^n, \quad m = p+1, \dots, q.$$

Passing to the limit and using Lemma 4.2, we see that  $r$  is optimal for the continuous relaxed problem. If there are no state constraints, by taking a sequence converging to some continuous optimal control, we arrive directly to the same conclusion.

**Lemma 4.3** (Consistency of adjoints and functional derivatives) Under Assumptions 2.2-5 and 3.1, if  $r^n \rightarrow r$  in  $R$ , then the corresponding discrete adjoint states  $z_-^n, z_+^n, z_{1-\theta}^n, z_\wedge^n$  converge to  $z_r$  in  $L^2(I, L^4(\Omega))$  (resp.  $L^2(Q)$ ) strongly if  $b \neq 0$  (resp.  $b = 0$ ), and  $z_{1-\theta}^n \rightarrow z_r$  in  $L^2(I, V)$  strongly. If  $r^n \rightarrow r$  and  $r^m \rightarrow r'$ , then

$$(4.32) \quad \lim_{n \rightarrow \infty} DG_m^n(r^n, r^m - r^n) = DG_m(r, r' - r), \quad m = 0, \dots, q.$$

**Proof.** The proof is similar to that of Lemma 4.2, using also the consistency of states.

Next, we study the behavior in the limit of extremal discrete controls. Consider the discrete problem with state constraints, Case (b). We shall construct sequences of perturbations  $(\varepsilon_m^n)$  converging to zero and such that the discrete problem is feasible for every  $n$ . Let  $r^m \in R^n$  be any solution of the problem without state constraints

$$(4.33) \quad c^n := \min_{w^n \in W^n} \left\{ \sum_{m=1}^p [G_m^n(r^n)]^2 + \sum_{m=p+1}^q [\max(0, G_m^n(r^n))]^2 \right\},$$

and set

$$(4.34) \quad \varepsilon_m^n := G_m^n(r^m), \quad m = 1, \dots, p, \quad \varepsilon_m^n := \max(0, G_m^n(r^m)), \quad m = p+1, \dots, q.$$

Let  $\tilde{r}$  be an admissible control for the continuous relaxed problem, and  $(\tilde{r}^n \in R^n)$  a sequence converging to  $\tilde{r}$  in  $R$  (Proposition 4.1). We have

$$(4.35) \quad \lim_{n \rightarrow \infty} [G_m^n(\tilde{r}^n)]^2 = [G_m(\tilde{r})]^2 = 0, \quad m = 1, \dots, p,$$

$$(4.36) \quad \lim_{n \rightarrow \infty} [\max(0, G_m^n(\tilde{r}^n))]^2 = [\max(0, G_m(\tilde{r}))]^2 = 0, \quad m = p+1, \dots, q,$$

which imply a fortiori that  $c^n \rightarrow 0$ , hence  $\varepsilon_m^n \rightarrow 0$ ,  $m = 1, \dots, q$ . Then clearly the discrete problem, Case (b), is feasible for every  $n$ , for these perturbations  $\varepsilon_m^n$ . We suppose in the sequel that the perturbations  $\varepsilon_m^n$  are chosen as in the above minimum feasibility procedure. Note that in practice we usually have  $c^n = 0$ , for sufficiently large  $n$ , due to sufficient discrete controllability, in which case the perturbations  $\varepsilon_m^n$  are equal to zero, i.e. the discrete problem with zero perturbations is feasible.

**Theorem 4.2** Under Assumptions 2.2-5 and 3.1, for each  $n$ , let  $r^n$  be admissible and extremal for the discrete problem, Case (b). Then every relaxed accumulation point of  $(r^n)$  is admissible and extremal for the continuous relaxed problem.

**Proof.** Since  $R$  is compact and  $\sum_{m=0}^q |\lambda_m^n| = 1$ , let  $(r^n)$ ,  $(\lambda_m^n)$ ,  $m = 0, \dots, q$ , be subsequences such that  $r^n \rightarrow r$  in  $R$  and  $\lambda_m^n \rightarrow \lambda_m$ ,  $m = 0, \dots, q$ , and consider the discrete principle in global form, which can be written as

$$(4.37) \quad \int_Q H^n(x^n, t_\theta^n, y_\theta^n, \nabla y_\theta^n, z_{1-\theta}^n, r^m - r^n) dx dt \geq 0, \quad \forall r^m \in R^n.$$

Passing to the limit, by Lemmas 4.2, 4.3 and Proposition 2.1 in [4], we obtain

$$(4.38) \quad \int_Q H(x, t, y, \nabla y, z, r^l(x, t) - r(x, t)) dx dt \geq 0, \quad \forall r^m \in R^n.$$

On the other hand, we have similarly

$$(4.39) \quad \lambda_m G_m(r) = \lim_{n \rightarrow \infty} \lambda_m^n [G_m^n(r^n) - \varepsilon_m^n] = 0, \quad m = p+1, \dots, q,$$

$$(4.40) \quad G_m(r) = \lim_{n \rightarrow \infty} [G_m^n(r^n) - \varepsilon_m^n] = 0, \quad m = 1, \dots, p,$$

$$(4.41) \quad G_m(r) = \lim_{n \rightarrow \infty} [G_m^n(r^n) - \varepsilon_m^n] \leq 0, \quad m = p+1, \dots, q,$$

and  $\lambda_0 \geq 0$ ,  $\lambda_m \geq 0$ ,  $m = p+1, \dots, q$ ,  $\sum_{m=0}^q |\lambda_m| = 1$ , which show that  $r$  is admissible and extremal for the continuous relaxed problem.

## 5 Discrete penalized conditional descent methods

Let  $(M_m^l)$ ,  $m = 1, \dots, q$ , be positive increasing sequences such that  $M_m^l \rightarrow \infty$  as  $l \rightarrow \infty$ , and define the *penalized discrete functionals*

$$(5.1) \quad G^{nl}(r^n) := G_0^n(r^n) + \left\{ \sum_{m=1}^p M_m^l [G_m^n(r^n)]^2 + \sum_{m=p+1}^q M_m^l [\max(0, G_m^n(r^n))]^2 \right\} / 2.$$

Let  $b', c' \in (0, 1)$ , and let  $(\beta^l)$ ,  $(\zeta_k)$  be positive sequences, with  $(\beta^l)$  decreasing and converging to zero, and  $\zeta_k \leq 1$ . The algorithm described below contains two versions. In the case of the progressively refining version, we shall make the following assumptions.

**Assumptions 5.1** The (possibly) finer discretization for  $n+1$  is defined by subdividing the elements  $E_i^n$  into subelements (e.g. triangles into 4 triangles) and by slightly (up to  $o(h^n)$ , as  $\Gamma$  is  $C^1$  or piecewise  $C^1$ ) transforming the resulting boundary elements so as to fit  $\Gamma^{n+1}$ , and then by setting  $N^{n+1} = \kappa N^n$ , for some integer  $\kappa \geq 2$ . The *discrete penalized conditional descent methods* are described in the following algorithm.

**Algorithm**

*Step 1.* Set  $k := 0$ ,  $l := 1$ , choose a value of  $n$  and an initial control  $r_0^n \in R^n$ .

*Step 2.* Find  $r_k^{nl} \in R^n$  such that

$$(5.2) \quad d_k := DG^{nl}(r_k^{nl}, \bar{r}_k^{nl} - r_k^{nl}) = \min_{r^m \in R^n} DG^{nl}(r_k^{nl}, r^m - r_k^{nl}).$$

*Step 3.* If  $|d_k| \leq \beta^l$ , set  $r^{nl} := r_k^{nl}$ ,  $\bar{r}^{nl} := \bar{r}_k^{nl}$ ,  $d^l := d_k$ ; [if the discretization for  $n+1$  is finer, set first  $\tilde{r}_k := r_k^{nl}$  on  $\Omega$ , and then define  $r_k^{n+1, l+1}$  as the modified control resulting from  $\tilde{r}_k$  after the slight transformation in the construction of the new boundary elements  $E_i^{n+1}$ ; set  $n := n+1$ ]; set  $l := l+1$  and go to Step 2. If  $|d_k| > \beta^l$ , go to Step 4.

*Step 4. (Armijo step search)* Find the lowest integer value  $s \in \mathbb{Z}$ , say  $\bar{s}$ , such that  $\alpha(s) := c^{s\bar{s}} \zeta_k \in (0, 1]$  and  $\alpha(s)$  satisfies the inequality

$$(5.3) \quad G^n(w_k^n + \alpha(s)(v_k^n - w_k^n)) - G^n(w_k^n) \leq \alpha(s)b'd_k,$$

and then set  $\alpha_k := \alpha(\bar{s})$ .

*Step 5.* Choose any  $r_{k+1}^{nl} \in R^n$  such that

$$(5.4) \quad G^{nl}(r_{k+1}^{nl}) \leq G^{nl}(r_k^{nl} + \alpha(s)(\bar{r}_k^{nl} - r_k^{nl})),$$

set  $k := k+1$ , and go to Step 2.

In the above Algorithm, we consider two versions:

**Version A.** [ $n := n+1$  etc.] is *skipped* in Step 3:  $n$  is a constant integer chosen in Step 1, i.e. we choose a *fixed discretization* and replace the discrete functionals  $G_m^n$  by the perturbed ones  $\tilde{G}_m^n := G_m^n - \varepsilon_m^n$ .

**Version B.** [ $n := n+1$  etc.] is *not skipped* in Step 3: we have a *progressively refining* discrete method, i.e.  $n \rightarrow \infty$  (see proof of Theorem 5.1 below), in which case we can take  $n=1$  in Step 1, hence  $n=l$  in the Algorithm. This version has the advantage of reducing computing time and memory, and also of avoiding the computation of the minimum feasibility perturbations  $\varepsilon_m^n$ . It is justified by the fact that finer discretizations become more efficient as the iterate gets closer to an extremal control, while coarser ones in the early iterations have not much influence on the final results.

One can easily see that a *classical* control  $\bar{r}_k^{nl}$  in Step 2 can be found for every  $k$  by minimizing in  $u \in U$  the numerical integral on  $E_i^n$

$$(5.5) \quad \mu(E_i^n) \sum_{v=1}^s C^v H(t_{j\theta}^n, x_i^v, y_{j\theta}^{nl}(x_i^{vn}), \nabla y_{j\theta}^{nl}(x_i^{vn}), u)$$

independently for each  $i=1, \dots, M$ ,  $j=1, \dots, N$ . On the other hand, since clearly  $d_k \leq 0$  and  $b' \in (0, 1)$ , by the definition of the directional derivative the Armijo step  $\alpha_k$  in Step 4 can be found for every  $k$ , if  $d_k \neq 0$ .

A continuous or discrete extremal control is called *abnormal* if there exist multipliers as in the corresponding optimality conditions, with  $\lambda_0 = 0$  (or  $\lambda_0^n = 0$ ). A control is admissible *and* abnormal extremal in exceptional, degenerate, situations (see [20]).

With  $w^{nl}$  defined in Step 3, define the *sequences of multipliers*

$$(5.6) \quad \lambda_m^{nl} := M_m^l G_m^n(r^{nl}), \quad m = 1, \dots, p, \quad \lambda_m^{nl} := M_m^l \max(0, G_m^n(r^{nl})), \quad m = p+1, \dots, q,$$

**Theorem 5.1** We suppose that Assumptions 2.2-5, 3.1, and 5.1 (progressively refining case) are satisfied.

(i) In Version B, let  $(r^{nl})$  be a subsequence, considered as a sequence in  $R$ , of the sequence generated by the Algorithm in Step 3 that converges to some  $r$  in  $R$ , as  $l \rightarrow \infty$  (hence  $n \rightarrow \infty$ ). If the sequences  $(\lambda_m^{nl})$  are bounded, then  $r$  is admissible and extremal for the continuous relaxed problem.

(ii) In Version A, let  $(r^{nl})$ ,  $n$  fixed, be a subsequence of the sequence generated by the Algorithm in Step 3 that converges to some  $r^n \in R^n$  as  $l \rightarrow \infty$ . If the sequences  $(\lambda_m^{nl})$  are bounded, then  $r^n$  is admissible and extremal for the fixed discrete problem.

(iii) In any of the two above convergence cases (i), (ii), suppose that the (discrete or continuous) limit problem has no admissible, abnormal extremal, controls. If the limit control is admissible, then the sequences of multipliers are bounded, and this control is extremal as above.

**Proof.** We shall first show that  $l \rightarrow \infty$  in the Algorithm. Suppose, on the contrary, that  $l$ , hence  $n$  (in both Versions A, B), remains constant after a finite number of iterations in  $k$ , and so we drop here the indices  $l$  and  $n$ . Let us show that then  $d_k \rightarrow 0$ . Since  $R$  is compact, let  $(r_k)_{k \in K}$ ,  $(\bar{r}_k)_{k \in K}$  be subsequences of the sequences generated in Steps 2 and 5 such that  $r_k \rightarrow \tilde{r}$ ,  $\bar{r}_k \rightarrow \tilde{\bar{r}}$ , in  $R$ , as  $k \rightarrow \infty$ ,  $k \in K$ . Clearly, by Step 2,  $d_k \leq 0$  for every  $k$ , hence

$$(5.7) \quad d := \lim_{k \rightarrow \infty, k \in K} d_k = DG(\tilde{r}, \tilde{\bar{r}} - \tilde{r}) \leq 0.$$

Suppose that  $d < 0$ . The function  $\Phi(\alpha) := G(r + \alpha(r' - r))$  is continuous on  $[0, 1]$ . Since the directional derivative  $DG(r, r' - r)$  is linear w.r.t.  $r' - r$ ,  $\Phi$  is differentiable on  $(0, 1)$  and has derivative  $\Phi'(\alpha) = DG(r + \alpha(r' - r), r' - r)$ . Using the mean value theorem, we have, for each  $\alpha \in (0, 1]$

$$(5.8) \quad G(r_k + \alpha(\bar{r}_k - w_k)) - G(r_k) = \alpha DG(r_k + \alpha'(\bar{r}_k - r_k), \bar{r}_k - r_k),$$

for some  $\alpha' \in (0, \alpha)$ . Therefore, for  $\alpha \in [0, 1]$ , by the continuity of  $DG$  (Lemma 3.1)

$$(5.9) \quad G(r_k + \alpha(\bar{r}_k - r_k)) - G(r_k) = \alpha(d + \varepsilon_{k\alpha}),$$

where  $\varepsilon_{k\alpha} \rightarrow 0$  as  $k \rightarrow \infty$ ,  $k \in K$ , and  $\alpha \rightarrow 0^+$ . Now, we have  $d_k = d + \eta_k$ , where  $\eta_k \rightarrow 0$  as  $k \rightarrow \infty$ ,  $k \in K$ , and since  $b' \in (0, 1)$

$$(5.10) \quad d + \varepsilon_{k\alpha} \leq b(d + \eta_k) = b'd_k,$$

for  $\alpha \in [0, \bar{\alpha}]$ , for some  $\bar{\alpha} > 0$ , and  $k \geq \bar{k}$ ,  $k \in K$ . Hence

$$(5.11) \quad G(r_k + \alpha(\bar{r}_k - r_k)) - G(r_k) \leq \alpha b'd_k,$$

for  $\alpha \in [0, \bar{\alpha}]$ , for some  $\bar{\alpha} > 0$ , and  $k \geq \bar{k}$ ,  $k \in K$ . It follows from the choice of the Armijo step  $\alpha_k$  in Step 4 that  $\alpha_k \geq c\bar{\alpha}$ , for  $k \geq \bar{k}$ ,  $k \in K$ . Hence

$$(5.12) \quad G(r_{k+1}) - G(r_k) \leq G(r_k + \alpha_k(\bar{r}_k - r_k)) - G(r_k) \leq \alpha_k b'd_k \leq c\bar{\alpha} b'd_k \leq c\bar{\alpha} b'd/2,$$



for  $k \geq \bar{k}$ ,  $k \in K$ . It follows that  $G(r_k) \rightarrow -\infty$  as  $k \rightarrow \infty$ ,  $k \in K$ , which contradicts the fact that  $G(r_k) \rightarrow G(\tilde{r})$  as  $k \rightarrow \infty$ ,  $k \in K$ , by the continuity of the discrete functional (Lemma 3.1). Therefore, we must have  $d = 0$ , and  $d_k \rightarrow d = 0$ , for the whole sequence, since the limit 0 is unique. But Step 3 then implies that  $l \rightarrow \infty$ , which is a contradiction. Therefore,  $l \rightarrow \infty$ . This shows also that  $n \rightarrow \infty$  in Version B.

(i) Let  $(r^{nl})$  be a subsequence (same notation) of the sequence generated in Step 3, that converges to some accumulation point  $r \in R$  as  $l, n \rightarrow \infty$ . Suppose that the sequences  $(\lambda_m^{nl})$  are bounded and (up to subsequences) that  $\lambda_m^{nl} \rightarrow \lambda_m$ . By Lemma 4.2, we have

$$(5.13) \quad 0 = \lim_{l \rightarrow \infty} \frac{\lambda_m^{nl}}{M_m^l} = \lim_{l \rightarrow \infty} G_m^n(r^{nl}) = G_m(r), \quad m = 1, \dots, p,$$

$$(5.14) \quad 0 = \lim_{l \rightarrow \infty} \frac{\lambda_m^{nl}}{M_m^l} = \lim_{l \rightarrow \infty} [\max(0, G_m^n(r^{nl}))] = \max(0, G_m(r)), \quad m = p+1, \dots, q,$$

which show that  $r$  is admissible. Now, by Steps 2 and 3 we have, for every  $v^m \in W^n$

$$(5.15) \quad \begin{aligned} DG^{nl}(r^{nl}, r^m - r^{nl}) \\ = DG_0^n(r^{nl}, r^m - r^{nl}) + \sum_{m=1}^p \lambda_m^{nl} DG_m^n(r^{nl}, r^m - r^{nl}) + \sum_{m=p+1}^q \lambda_m^{nl} DG_m^n(r^{nl}, r^m - r^{nl}) \\ \geq d^l. \end{aligned}$$

Using Lemmas 4.2, 4.3 and Proposition 2.1 in [4], we can pass to the limit in this inequality as  $l, n \rightarrow \infty$  and obtain

$$(5.16) \quad DG_0(r, r^1 - r) + \sum_{m=1}^p \lambda_m DG_m(r, r^1 - r) + \sum_{m=p+1}^q \lambda_m DG_m(r, r^1 - r) \geq 0.$$

By construction of the  $\lambda_m^{nl}$ , we clearly have in the limit  $\lambda_0 = 1$ ,  $\lambda_m \geq 0$ ,  $m = p+1, \dots, q$ ,  $\sum_{m=0}^q |\lambda_m| := c \geq 1$ , and we can suppose that  $\sum_{m=0}^q |\lambda_m| = 1$  by dividing the above inequality by  $c$ . On the other hand, if  $G_m(r) < 0$ , for some index  $m \in [p+1, q]$ , then for sufficiently large  $l$  we have  $G_m^{nl}(r^{nl}) < 0$  and  $\lambda_m^l = 0$ , hence  $\lambda_m = 0$ , i.e. the transversality conditions hold. Therefore,  $r$  is also extremal.

(ii) The admissibility of the limit control  $r^n$  is proved as in (i). Passing here to the limit in the inequality resulting from Step 2 as  $l \rightarrow \infty$ , for  $n$  fixed, and using Lemmas 3.1 and 3.2, we obtain, similarly to (i) (with  $\lambda_0 = 1$ )

$$(5.17) \quad \sum_{m=0}^q \lambda_m D\tilde{G}_m^n(r^n, r^m - r^n) = \sum_{m=0}^q \lambda_m DG_m^n(r^n, r^m - r^n) \geq 0, \quad \forall r^m \in R^n,$$

with multipliers as in the optimality conditions, and the discrete transversality conditions

$$(5.18) \quad \lambda_m^n \tilde{G}_m^n(r^n) = \lambda_m^n [G_m^n(r^n) - \varepsilon_m^n] = 0, \quad m = p+1, \dots, q,$$

(iii) In either of the above convergence cases (i) or (ii), suppose that the limit control is admissible and that the limit problem has no admissible, abnormal extremal, controls. Suppose that the multipliers are not all bounded. Then, dividing the corresponding inequality resulting from Step 2 by the greatest multiplier norm and passing to the limit for a subsequence, we see that we obtain an optimality inequality

where the first multiplier is zero, and that the limit control is abnormal extremal, a contradiction. Therefore, the sequences of multipliers are bounded, and by (i) or (ii), this limit control is extremal as above.

In practice, by choosing moderately growing sequences  $(M_m^l)$  and a sequence  $(\beta^l)$  relatively fast converging to zero, the resulting sequences of multipliers  $(\lambda_m^{nl})$  are often kept bounded. One can choose a fixed  $\zeta_k := \zeta \in (0,1]$  in Step 4; a usually faster and adaptive procedure is to set  $\zeta_0 := 1$  and  $\zeta_k := \alpha_{k-1}$ , for  $k \geq 1$ .

The Algorithm can be implemented as follows. Suppose that the integrals on  $\Omega$  involved in the discrete state equation and the functionals are calculated with sufficient accuracy by using an integration rule of the form

$$(5.19) \quad \int_{\Omega} \phi(x) dx \approx \sum_{i=1}^M [\text{meas}(E_i^n) \sum_{v=1}^s C^v \phi(x_i^{vn})].$$

We first choose the initial discrete control in Step 1 to be of Gamkrelidze type, i.e. equal on each block  $Q_{ij}^n$  to a convex combination of  $(s+q+1)+1$  ( $s$  integration nodes) Dirac measures on  $U$  concentrated at  $(s+q+1)+1$  points of  $U$ . Suppose, by induction, that the control  $r_k^{nl}$  computed in the Algorithm is of Gamkrelidze type. Since the control  $\bar{r}_k^{nl}$  in Step 2 is chosen to be classical, i.e. blockwise Dirac (see above), the control  $\tilde{r}_k^{nl} := (1-\alpha_k)r_k^{nl} + \alpha_k\bar{r}_k^{nl}$  in Step 5 is blockwise equal to a convex combination of  $(s+q+1)+2$  Dirac measures. Using now a known property of convex hulls of finite vector sets, we can construct a Gamkrelidze control  $r_{k+1}^{nl}$  equivalent to  $\tilde{r}_k^{nl}$ , i.e.  $\tilde{r}_k^{nl}$ , i.e. such that the following  $s+q+1$  equalities (i.e. equality in  $\mathbb{R}^{s+q+1}$ ) hold

$$(5.20) \quad f(t_{\theta}^n, x_i^{vn}, \tilde{y}_{\theta}^{nl}(x_i^{vn}), r_{k+1,i}^{nl}) = f(t_{\theta}^n, x_i^{vn}, \tilde{y}_{\theta}^{nl}(x_i^{vn}), \tilde{r}_{ki}^{nl}), \quad v=1, \dots, s,$$

$$(5.21) \quad \mu(E_i^n) \sum_{v=1}^s C^v g_m(t_{\theta}^n, x_i^{vn}, \tilde{y}_{\theta}^{nl}(x_i^{vn}), \nabla \tilde{y}_{\theta}^{nl}(x_i^{vn}), r_{k+1,i}^{nl}) \\ = \mu(E_i^n) \sum_{v=1}^s C^v g_m(t_{\theta}^n, x_i^{vn}, \tilde{y}_{\theta}^{nl}(x_i^{vn}), \nabla \tilde{y}_{\theta}^{nl}(x_i^{vn}), \tilde{r}_{ki}^{nl}), \quad m=0, \dots, q,$$

for each  $i=1, \dots, M$ ,  $j=1, \dots, N$ , where  $\tilde{y}^{nl}$  corresponds to  $\tilde{r}_k^{nl}$ , by selecting only  $(s+q+1)+1$  appropriate points in  $U$  among the  $(s+q+1)+2$  ones defining  $\tilde{r}_k^{nl}$ . Then the control  $r_{k+1}^{nl}$  clearly yields the same discrete state and functionals as  $\tilde{r}_k^{nl}$  and thus satisfies Step 5. Therefore, the constructed control  $r_k^{nl}$  is of Gamkrelidze type for every  $k$  (note also that by the construction of the control  $r_k^{n+1, l+1}$  in Step 3 if the discretization is refined, this control is still of Gamkrelidze type, but w.r.t. to the new elements  $E_i^{n+1}$ ). Finally, discrete Gamkrelidze controls computed as above can then be approximated by subblockwise (w.r.t.  $t$ ) constant classical controls using a standard procedure (see [8]).

## 6 Numerical examples

Let  $\Omega := I := (0,1)$  and consider the following examples.

*Example 1.* Define the reference control and state

$$(6.1) \quad \bar{w}(x) := \begin{cases} 1, & \text{if } 0 \leq t \leq 0.5, \\ 1 - 2(t - 0.5)(0.2x + 0.4), & \text{if } 0.5 < t \leq 1, \end{cases} \quad \bar{y}(x) := x(1-x)e^{-t},$$

and consider the following optimal control problem, with state equation

$$(6.2) \quad y_t - y_{xx} + 0.5y|y| + (1+w-\bar{w})y \\ = 0.5\bar{y}|\bar{y}| + \bar{y} + [-x(1-x) + 2]e^{-t} + \sin y - \sin \bar{y} + 3(w-\bar{w}),$$

$$(6.3) \quad y(x,t) = 0 \quad \text{on } \Sigma, \quad y(0,x) = x(1-x) \quad \text{in } \Omega,$$

nonconvex control constraint set

$$(6.4) \quad U := [0, 0.25] \cup [0.75, 1] \quad (\text{or } U := \{0, 1\}, \text{ on/off type control}),$$

and nonconvex cost functional to be minimized

$$(6.5) \quad G_0(u) := \int_Q \{0.5[(y-\bar{y})^2 + |\nabla y - \nabla \bar{y}|^2] - (w-0.5)^2 + 0.25\} dxdt$$

One can easily verify that the unique optimal *relaxed* control  $r$  is given by

$$(6.6) \quad r(x,t)\{1\} := \bar{w}(x,t), \quad r(x,t)\{0\} := 1 - r(x,t)\{1\},$$

for  $x \in Q$ , with optimal state  $\bar{y}$  and cost 0, and we see that  $r$  is concentrated at the two points 1 and 0;  $r$  is classical ( $\equiv 1$ ) if  $0 \leq t \leq 0.5$ , and non-classical otherwise. Note also that the optimal cost value 0 can be approximated as closely as possible by using a classical control, as  $W$  is dense in  $R$ , but cannot be attained for such a control because the control values  $u \in (0.25, 0.75)$ , or  $u \in (0, 1)$ , (of  $\bar{w}$ ) do not belong to  $U$ .

The Algorithm, without penalties, was applied to this problem using the midpoint integration rule on each interval  $E_i^n \subset \Omega$ , with step sizes  $h = \Delta t = 1/100$ ,  $\theta$ -scheme parameter  $\theta = 1$  (implicit Euler method), Armijo parameters  $b' = c' = 0.5$ , and constant initial control  $r_0^n(x,t) := 0.5(\delta_0 + \delta_1)$ ,  $(x,t) \in Q$ , where  $\delta_0, \delta_1$  are the Dirac measures at 0 and 1. After 90 iterations in  $k$ , we obtained the results:

$$(6.7) \quad G_0^n(r_k^n) = 3.471 \cdot 10^{-5}, \quad d_k = -1.480 \cdot 10^{-4}, \quad \eta_k = 5.722 \cdot 10^{-3},$$

where  $d_k$  was defined in the Algorithm and  $\eta_k$  is the discrete max state error at the points  $(ih, j\Delta t)$ . Figure 1 shows the last control probability function, for  $x = 0.5$  (cross-section), and we have  $p_0(x,t) := r_k^n(x,t)\{1\} = 1 - p_1(x,t)$ ; the other cross-sections are similar. Figure 2 shows the last computed state ( $\approx \bar{y}$ ).

*Example 2.* We introduce the equality state constraint

$$(6.8) \quad G_1(w) := \int_Q y dxdt = 0,$$

in Example 1. The continuous relaxed problem is feasible, as  $G_1(r_1) \approx 0.125 > 0$ ,  $G_1(r_0) \approx -0.103 < 0$ , where  $r_1(x,t) := \delta_1$ ,  $r_0(x,t) := \delta_0$ ,  $(x,t) \in Q$ , and the function  $\phi(\lambda) := G_1(\lambda r_1 + (1-\lambda)r_0)$  is continuous on  $[0,1]$ . Applying here the penalized Algorithm and the same parameters as in Example 1, we obtained after 147 iterations in  $k$  the results:

$$(6.9) \quad G_0^n(r_k^{nl}) = 7.526070032039354 \cdot 10^{-2}, \quad G_1^n(r_k^{nl}) = 6.667 \cdot 10^{-5}, \\ d_k = -2.058 \cdot 10^{-3}.$$

Figure 3 shows the last control probability function  $p_1(x,t) := r_k^{nl}(x,t)\{1\}$ , for  $x = 0.5$ , and we have also  $p_0(x,t) := r_k^{nl}(x,t)\{1\} = 1 - p_1(x,t)$ . Figure 4 shows the last computed state.

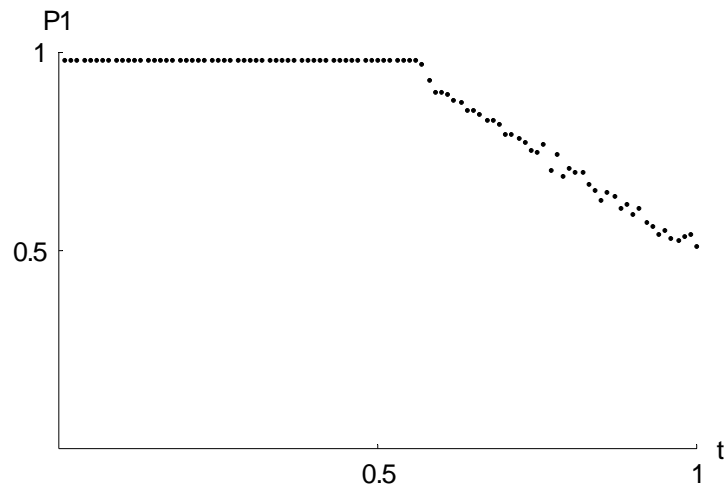


Figure 1. Example 1: Last relaxed control probability  $p_1$ , for  $x = 0.5$

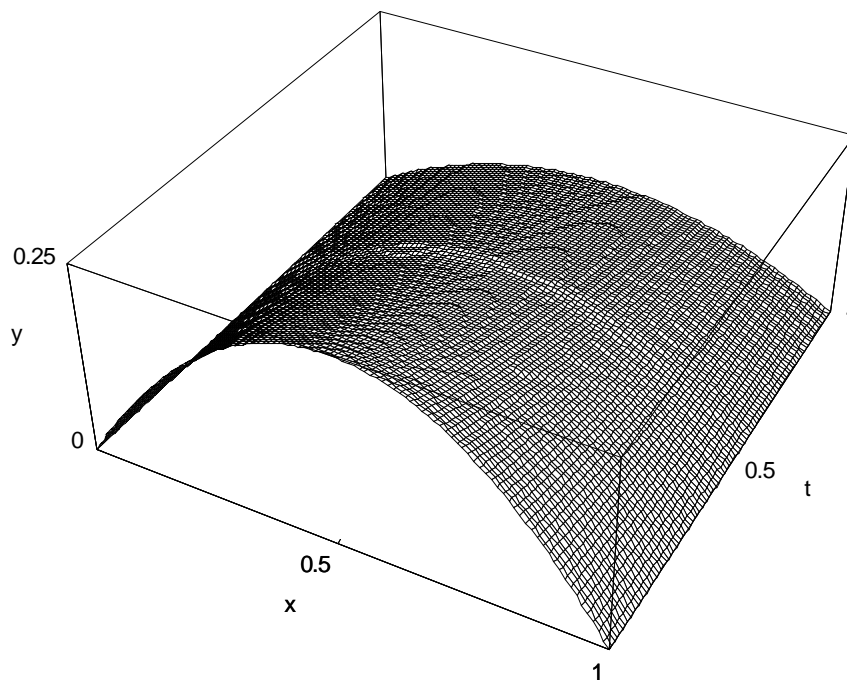


Figure 2. Example 1: Last state.

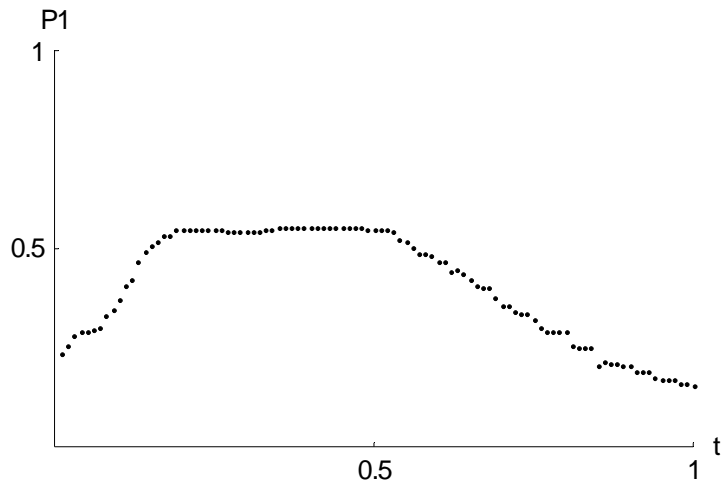


Figure 3. Example 2: Last relaxed control probability  $p_1$ , for  $x = 0.5$

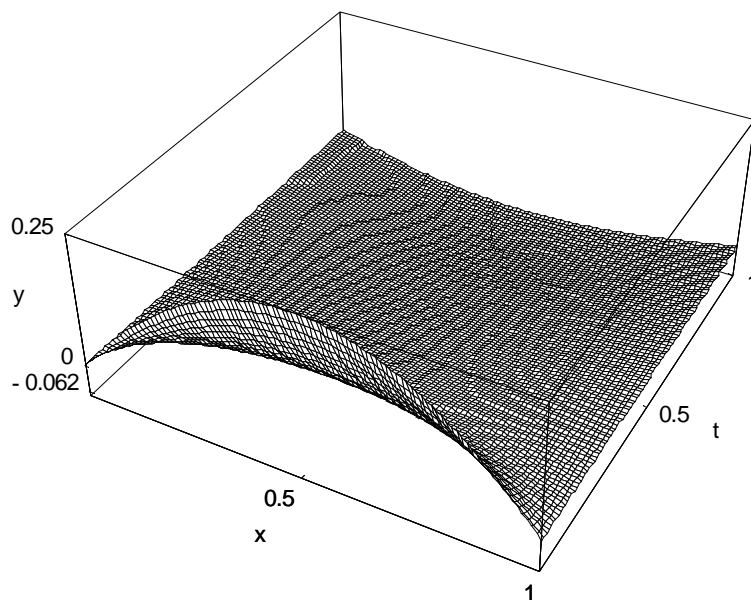


Figure 4. Example 2: Last state.

Finally, the progressively refining version of the algorithm was also applied to the above problems, with successive step sizes  $h = \Delta t = 1/25, 1/50, 1/100$ , in three equal iteration periods, and yielded results of similar accuracy, but required here less than half the computing time.

## References

- [1] R. A. Adams, Sobolev Spaces, Academic Press, New York, 1975.
- [2] S. Bartels, Adaptive approximation of Young measure solution in scalar non-convex variational problems, *SIAM J. Numer. Anal.*, 42 (2004) 505-629.
- [3] C. Cartensen and T. Roubíček, Numerical approximation of Young measure in nonconvex variational problems, *Numer. Math.*, 84 (2000) 395-415.
- [4] I. Chrysoverghi, Nonconvex optimal control of nonlinear monotone parabolic systems, *Systems Control Lett.*, 8 (1986) 55-62.
- [5] I. Chrysoverghi and A. Bacopoulos, Approximation of relaxed nonlinear parabolic optimal control problems, *J. Optim. Theory Appl.*, 77, 1 (1993) 31-50.
- [6] I. Chrysoverghi, A. Bacopoulos, B. Kokkinis and J. Coletsos, Mixed Frank-Wolfe penalty method with applications to nonconvex optimal control problems, *J. Optim. Theory Appl.*, 94, 2 (1997) 311-334.
- [7] I. Chrysoverghi, A. Bacopoulos, J. Coletsos and B. Kokkinis, Discrete approximation of nonconvex hyperbolic optimal control problems with state constraints, *Control Cybernet.*, 27, 1 (1998) 29-50.
- [8] I. Chrysoverghi, J. Coletsos and B. Kokkinis, Discrete relaxed method for semilinear parabolic optimal control problems, *Control Cybernet.*, 28, 2 (1999) 157-176.
- [9] I. Chrysoverghi, Discretization methods for semilinear parabolic optimal control problems, to appear in *Int. J. Numer. Anal. Modeling* (2005).
- [10] P. G. Ciarlet, *The Finite Element Method for Elliptic Problems*, North-Holland, New York, 1978.
- [11] H.O. Fattorini, *Infinite dimensional optimization theory and Optimal control*, Cambridge Univ. Press, Cambridge, 1999.
- [12] J. Mach, Numerical solution of a class of nonconvex variational problems by SQP, *Numer. Funct. Anal. Optim.*, 23 (2002) 573-587.
- [13] A.-M. Mataché, T. Roubíček and C. Schwab, Higher-order convex approximations of Young measures in optimal control, *Adv. Comput. Math.*, 19 (2003) 73-79.
- [14] R.A. Nicolaidis and N.J. Walkington, Strong convergence of numerical solutions to degenerate variational problems, *Math. Comp.* 64 (1995) 117-127.
- [15] E. Polak, *Optimization: Algorithms and Consistent Approximations*, Springer, Berlin, 1997.
- [16] T.R. Rockafellar and R. Wetts, *Variational Analysis*, Springer, Berlin, 1998.
- [17] T. Roubíček, *Relaxation in Optimization Theory and Variational Calculus*, Walter de Gruyter, Berlin, 1997.
- [18] R. Temam, *Navier-Stokes Equations*, North-Holland, New York, 1977.
- [19] V. Thomee, *Galerkin Finite Element Methods for Parabolic Problems*, Springer, Berlin, 1997.
- [20] J. Warga, *Optimal Control of Differential and Functional Equations*, Academic Press, New York, 1972.
- [21] J. Warga, Steepest descent with relaxed controls, *SIAM J. Control Optim.* 15 (1977) 674-682.