

## THE SWITCH POINT ALGORITHM\*

MAHYA AGHAE<sup>†</sup> AND WILLIAM W. HAGER<sup>†</sup>

**Abstract.** The switch point algorithm is a new approach for solving optimal control problems whose solutions are either singular or bang-bang or both singular and bang-bang, and which possess a finite number of jump discontinuities in an optimal control at the points in time where the solution structure changes. Problems in this class can often be reduced to an optimization over the switching points. Formulas are derived for the derivative of the objective with respect to the switch points, the initial costate, and the terminal time. All these derivatives can be computed simultaneously in just one integration of the state and costate dynamics. Hence, gradient-based unconstrained optimization techniques, including the conjugate gradient method or quasi-Newton methods, can be used to compute an optimal control. The performance of the algorithm is illustrated using test problems with known solutions and comparisons with other algorithms from the literature.

**Key words.** switch point algorithm, singular control, bang-bang control, total variation regularization

**AMS subject classifications.** 49M25, 49M37, 65K05, 90C30

**DOI.** 10.1137/21M1393315

**1. Introduction.** Let us consider a fixed terminal time control problem of the form

$$(1.1) \quad \min C(\mathbf{x}(T)) \quad \text{subject to} \quad \dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t)), \quad \mathbf{x}(0) = \mathbf{x}_0, \quad \mathbf{u}(t) \in \mathcal{U}(t),$$

where  $\mathbf{x} : [0, T] \rightarrow \mathbb{R}^n$  is absolutely continuous,  $\mathbf{u} : [0, T] \rightarrow \mathbb{R}^m$  is essentially bounded,  $C : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\mathbf{f} : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ , and  $\mathcal{U}(t)$  is a closed and bounded set for each  $t \in [0, T]$ . The dynamics  $\mathbf{f}$  and the objective  $C$  are assumed to be differentiable. The costate equation associated with (1.1) is the linear differential equation

$$(1.2) \quad \dot{\mathbf{p}}(t) = -\mathbf{p}(t) \nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t)), \quad \mathbf{p}(T) = \nabla C(\mathbf{x}(T)),$$

where  $\mathbf{p} : [0, T] \rightarrow \mathbb{R}^n$  is a row vector, the objective gradient  $\nabla C$  is a row vector, and  $\nabla_{\mathbf{x}} \mathbf{f}$  denotes the Jacobian of the dynamics with respect to  $\mathbf{x}$ . At the end of the introduction, the notation and terminology are summarized. Under the assumptions of the Pontryagin minimum principle, a local minimizer of (1.1) and the associated costate have the property that

$$H(\mathbf{x}(t), \mathbf{u}(t), \mathbf{p}(t)) = \inf\{H(\mathbf{x}(t), \mathbf{v}, \mathbf{p}(t)) : \mathbf{v} \in \mathcal{U}(t)\}$$

for almost every  $t \in [0, T]$ , where  $H(\mathbf{x}, \mathbf{u}, \mathbf{p}) = \mathbf{p}\mathbf{f}(\mathbf{x}, \mathbf{u})$  is the Hamiltonian.

The switch point algorithm is well suited for problems where an optimal control is piecewise smooth with one or more jump discontinuities at a finite set of times  $0 < s_1 < s_2 < \cdots < s_k < T$  where there is a fundamental change in the solution

\*Received by the editors January 21, 2021; accepted for publication (in revised form) May 2, 2021; published electronically July 13, 2021.

<https://doi.org/10.1137/21M1393315>

**Funding:** This research has been partially supported by National Science Foundation research grants 1819002 and 2031213.

<sup>†</sup>Department of Mathematics, University of Florida, Gainesville, FL 32611-8105 USA (mahyaaghaee@ufl.edu, <https://people.clas.ufl.edu/mahyaaghaee/>; hager@ufl.edu, <http://people.clas.ufl.edu/hager/>).

structure. We also define  $s_0 = 0$  and  $s_{k+1} = T$ . For illustration, suppose that  $\mathcal{U}(t) = \{\mathbf{v} \in \mathbb{R}^m : \boldsymbol{\alpha}(t) \leq \mathbf{v} \leq \boldsymbol{\beta}(t)\}$ , where  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta} : [0, T] \rightarrow \mathbb{R}^m$ , and  $\mathbf{f}(\mathbf{x}, \mathbf{u}) = \mathbf{g}(\mathbf{x}) + \mathbf{B}(\mathbf{x})\mathbf{u}$  with  $\mathbf{B} : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times m}$  and  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ . The switching function, the coefficient of  $\mathbf{u}$  in the Hamiltonian, is given by  $\mathcal{S}(t) = \mathbf{p}(t)\mathbf{B}(\mathbf{x}(t))$ . Under the assumptions of the Pontryagin minimum principle, an optimal solution of (1.1) has the property that

$$\begin{aligned} u_i(t) &= \alpha_i(t) \text{ if } \mathcal{S}_i(t) > 0, \\ u_i(t) &= \beta_i(t) \text{ if } \mathcal{S}_i(t) < 0, \text{ and} \\ u_i(t) &\in [\alpha_i(t), \beta_i(t)] \text{ otherwise,} \end{aligned}$$

for each index  $i \in [1, m]$  and for almost every  $t \in [0, T]$ . On intervals where  $\mathcal{S}_i$  is either strictly positive or strictly negative,  $u_i$  is uniquely determined from the minimum principle, and the control is said to be bang-bang. If for some  $i$ ,  $\mathcal{S}_i$  vanishes on an interval  $[\sigma, \tau]$ , then the problem is said to be singular, and on this singular interval, the first-order optimality conditions provide no information concerning  $u_i$  except that it satisfies the control constraints.

On any singular interval  $(s_j, s_{j+1})$ , not only does a component  $\mathcal{S}_i$  of the switching function vanish, but also [46] the derivatives of  $\mathcal{S}_i$  vanish, assuming they exist. If the singularity has finite order, then after equating derivatives of the switching function to zero, we eventually obtain a relationship of the form  $u_i(t) = \phi_{ij}(\mathbf{x}(t), \mathbf{p}(t), t)$  for all  $t \in (s_j, s_{j+1})$ . In vector notation, this relation can be expressed  $\mathbf{u}(t) = \boldsymbol{\phi}_j(\mathbf{x}(t), \mathbf{p}(t), t)$ . In many cases, it is possible to further simplify this to  $\mathbf{u}(t) = \boldsymbol{\phi}_j(\mathbf{x}(t), t)$ , where there is no dependence of the control on  $\mathbf{p}$ .

In the switch point algorithm, two separate cases are considered:

Case 1. For every  $j$ ,  $\boldsymbol{\phi}_j$  is independent of  $\mathbf{p}$ .

Case 2. For some  $j$ ,  $\boldsymbol{\phi}_j$  depends on  $\mathbf{p}$ .

In a nutshell, the switch point algorithm is based on the following observations. In Case 1, the control has the feedback form  $\mathbf{u}(t) = \boldsymbol{\phi}_j(\mathbf{x}(t), t)$  for  $t \in (s_j, s_{j+1})$ ,  $0 \leq j \leq k$ . For any given choice of the switching points, the solution of the state dynamics (assuming a solution exists) yields a value for the objective  $C(\mathbf{x}(T))$ . Hence, we can also think of the objective as a function  $C(\mathbf{s})$  depending on  $\mathbf{s}$ . In Case 2, where the control also depends on  $\mathbf{p}$ , we could (assuming a solution exists) integrate forward in time the coupled state and costate dynamics from any given initial condition  $\mathbf{p}(0) = \mathbf{p}_0$  with the control given by  $\mathbf{u}(t) = \boldsymbol{\phi}_j(\mathbf{x}(t), \mathbf{p}(t), t)$  on  $(s_j, s_{j+1})$ ,  $0 \leq j \leq k$ , to obtain a value for the objective. The objective would then be denoted  $C(\mathbf{s}, \mathbf{p}_0)$  since  $\mathbf{x}(T)$  depends on both  $\mathbf{s}$  and  $\mathbf{p}_0$ .

Suppose that  $\mathbf{s} = \mathbf{s}^*$  corresponds to the switching points for a solution of (1.1) and  $\mathbf{p}(0) = \mathbf{p}_0^*$  is the associated initial costate. In many applications, one finds that  $\mathbf{u}(t) = \boldsymbol{\phi}_j(\mathbf{x}(t), t)$  or  $\mathbf{u}(t) = \boldsymbol{\phi}_j(\mathbf{x}(t), \mathbf{p}(t), t)$  remain feasible in (1.1) for  $\mathbf{s}$  in a neighborhood of  $\mathbf{s}^*$  and for  $\mathbf{p}(0)$  in a neighborhood of  $\mathbf{p}_0^*$ . Moreover,  $C(\mathbf{s})$  or  $C(\mathbf{s}, \mathbf{p}_0)$  achieves a local minimum at  $\mathbf{s}^*$  or  $(\mathbf{s}^*, \mathbf{p}_0^*)$ . Therefore, at least locally, we could replace (1.1) by the problem of minimizing the objective over  $\mathbf{s}$  and  $\mathbf{p}_0$ .

We briefly review previous numerical approaches for singular control problems. Since the literature in the area is huge, our goal is to mostly highlight historical trends. One of the first approaches for singular control problems was what Jacobson, Gershwin, and Lele [30] called the  $\epsilon - \alpha(\cdot)$  algorithm, although today it would be called the proximal point method. The idea was to make a singular problem nonsingular by adding a strongly convex quadratic term to the objective. The new objective is

$$C(\mathbf{x}(T)) + \frac{\epsilon_k}{2} \int_0^T \|\mathbf{u}(t) - \mathbf{u}_k(t)\|^2 dt,$$

and the solution of the problem with the modified objective yields  $\mathbf{u}_{k+1}$ . Jacobson denoted the approximating control sequence by  $\alpha_k$ , so the scheme was referred to as the  $\epsilon - \alpha(\cdot)$  algorithm. The choice of  $\epsilon_k$  is a delicate issue; if it is too small, then  $\mathbf{u}_{k+1}$  can oscillate wildly, but if it is chosen just right, good approximations to solutions of singular control problems have been obtained.

In a different approach, Anderson [3] considers a problem where the control starts out nonsingular, and then changes to singular at switch time  $s_1$ . It is proposed to make a guess for the initial costate  $\mathbf{p}(0)$  and then adjust it until the junction conditions at the boundary between the singular and nonsingular control are satisfied. Then further adjustments to the initial costate are made to satisfy the terminal conditions on the costate. If these further adjustments cause the junction conditions to be violated unacceptably, then the entire process would be repeated, first modifying the initial costate to update the switching point and to satisfy the junction conditions, and then to further modify the initial costate to satisfy the terminal conditions. Aly [2] also considers the method of Anderson, but with more details. Maurer [37] uses multiple shooting techniques to satisfy junction conditions and boundary conditions in singular control problems. In general, choosing switching points and adjusting the initial costate condition to satisfy the junction and terminal conditions could be challenging.

Papers closer in spirit to the switch point algorithm include more recent works of Maurer et al. [38] and Vossen [50]. In both cases, the authors consider general boundary conditions of the form  $\phi(\mathbf{x}(0), \mathbf{x}(T), T) = \mathbf{0}$ , where  $\phi : \mathbb{R}^{2n+1} \rightarrow \mathbb{R}^r$ ,  $0 \leq r \leq 2n$ . In [38] the authors focus on bang-bang control problems, while Vossen considers singular problems where the control has the feedback form of Case 1. In both papers, the authors view the objective as a function of the switching points, and the optimization problem becomes a finite dimensional minimization over the switching points and the initial state subject to the boundary condition. Vossen in his dissertation [49] and in [51, Prop. 4.12] shows that when the switching points correspond to the switching points of a control for the continuous problem which satisfies the minimum principle, then the first-order optimality conditions are satisfied for the finite dimensional optimization problem. This provides a rigorous justification for the strategy of replacing the original infinite dimensional control problem by a finite dimensional optimization over the switching points and the boundary values of the state.

A fundamental difference between our approach and the approaches in [38] and [50] is that in the earlier work, the derivative of the objective is expressed in terms of the partial derivative of each state variables with respect to each switching point, where this matrix of partial derivatives is obtained by a forward propagation using the system dynamics. We circumvent the evaluation of the matrix of partial derivatives by using the costate equation to directly compute the partial derivative of the cost with respect to all the switching points; there is one forward integration of the state equation and one backward integration of the costate equation to obtain the partial derivatives of the objective with respect to all the switching points at the same time. In a sense, our approach is a generalization of [26, Thm. 2] which considers purely bang-bang controls. One benefit associated with the computation of the matrix of partial derivatives of each state with respect to each switch point is that with marginal additional work, second-order optimality conditions can be checked.

When the control has the feedback form of Case 1 and the dynamics are affine in the control, our formula for the objective derivative reduces to the product between the switching function and the jump in the control at the switching point.

Consequently, at optimality in the finite dimensional problem, the switching function should vanish. The vanishing of the switching function at the switching point is a classic necessary optimality condition [1, 39, 42]. Our formula, however, applies to arbitrary, not necessarily optimal controls, and hence it can be combined with a nonlinear programming algorithm to solve the control problem (1.1).

The methods discussed so far are known as indirect methods; they are indirect in the sense that steps of the algorithm employ information gleaned from the first-order optimality conditions. In a direct method, the original continuous-in-time problem is discretized using finite elements, finite differences, or collocation to obtain a finite dimensional problem which is solved by a nonlinear optimization code. These discrete problems can be difficult due to both ill conditioning and discontinuities in the optimal control at the switching points. If the location of the switching points were known, then high-order approximations are possible using orthogonal collocation techniques [25]. But in general, the switching points are not known, and mesh refinement techniques [17, 18, 32, 33, 40] can lead to highly refined meshes in a neighborhood of discontinuities, which can lead to a large dimension for the discrete problem in order to achieve a specified error tolerance.

Betts briefly touched on a hybrid direct/indirect approach to singular control in [9, sect. 4.14.1], where the switching points between singular and nonsingular regions are introduced as variables in the discrete problem, a junction condition is imposed at the switching points, and the form of the control in the singular region is explicitly imposed. A relatively accurate solution of the Goddard rocket problem [12] was obtained.

In a series of papers (see [4, 5, 14, 15] and the reference therein), Biegler and colleagues develop approaches to singular control problems that combine ideas from both direct and indirect schemes. In [14], the algorithm utilizes an inner problem where the mesh and the approximations to the switching points are fixed, and the discretized control problem is solved by a nonlinear optimization code such as IPOPT [52]. Then an outer problem is formulated where the finite element mesh is modified so as to reduce errors in the dynamics or make the switching function and its derivative closer to zero in the singular region. In [15] more sophisticated rules are developed for moving grid points or either inserting or deleting grid points by monitoring errors or checking for spikes. In [5], the inner and outer problems are also combined to form a single nonlinear program that is optimized. In [4] the direct method is mostly used to obtain a starting guess for the indirect phase of the algorithm. In the indirect phase, special structure is imposed on the control to reduce or eliminate wild oscillation, and conditions are imposed to encourage the vanishing of the switching function in the singular region and the constancy of the Hamiltonian.

In comparing the switch point algorithm to the existing literature, the necessary optimality conditions are not imposed except for our assumption, in the current version of the algorithm, that the control in the singular region has been expressed as a function of the state and/or costate. Note that the papers of Biegler and colleagues do not make this assumption; instead the vanishing of the switching function in the singular region is a constraint in their algorithms. The switch point algorithm focuses on minimizing the objective with respect to the switching points; presumably, the necessary optimality conditions will be satisfied at a minimizer of the objective, but these conditions are not imposed on the iterates.

After solving the switch point algorithm's finite dimensional optimization problem associated with a bang-bang or singular control problem, one may wish to check the optimality of the computed finite dimensional solution in the original continuous

control problem (1.1); to check whether a switching point was missed in the finite dimensional optimization problem, one could return to the original continuous control problem and test whether the minimum principle holds at other points in the domain, not just at the switching points, by performing an accurate integration of the differential equations between the switching points. To check the local optimality of the computed solution of the finite dimensional problem, one could test the second-order sufficient optimality conditions. As noted in [38, Thm. 3.1] for bang-bang control, satisfaction of the second-order sufficient optimality conditions in the finite dimensional problem implies that the associated solution of the continuous control problem is a strict minimum. Second-order sufficient optimality conditions for bang-singular controls are developed in [50].

The paper is organized as follows. In section 2, we explain how to compute the derivative of  $C$  with respect to  $\mathbf{s}$  (Case 1), while section 3 deals with the derivative of  $C$  with respect to both  $\mathbf{s}$  and  $\mathbf{p}_0$  (Case 2). Section 4 considers free terminal time problems and obtains the derivative of the objective with respect to the final time  $T$ . The efficient application of the switch point algorithm requires a good guess for the structure of an optimal control. In section 5 we explain one approach for generating a starting guess using total variation (TV) regularization, which has been effective in image reconstruction [45] at replacing blurry edges with sharp edges. Finally, section 6 provides some comparisons with other methods from the literature using test problems with known solutions.

**Notation and terminology.** By a valid choice of the switch points, we mean that  $\mathbf{s}$  satisfies the relations  $0 = s_0 < s_1 < s_2 < \cdots < s_k < s_{k+1} = T$ . Throughout the paper,  $\|\cdot\|$  is any norm on  $\mathbb{R}^n$ . The ball with center  $\mathbf{c} \in \mathbb{R}^n$  and radius  $\rho$  is given by  $\mathcal{B}_\rho(\mathbf{c}) = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x} - \mathbf{c}\| \leq \rho\}$ . The expression  $\mathcal{O}(\Delta s)$  denotes a quantity that is bounded in absolute value by  $c|\Delta s|$ , where  $c$  is a constant that is independent of  $\Delta s$ . Given  $\mathbf{x}$  and  $\mathbf{y} \in \mathbb{R}^n$ , we let  $[\mathbf{x}, \mathbf{y}]$  denote the line segment connecting  $\mathbf{x}$  and  $\mathbf{y}$ . In other words,

$$[\mathbf{x}, \mathbf{y}] = \{\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x}) : \alpha \in [0, 1]\}.$$

The Jacobian of  $\mathbf{f}(\mathbf{x}, \mathbf{u})$  with respect to  $\mathbf{x}$  is denoted  $\nabla_{\mathbf{x}}\mathbf{f}(\mathbf{x}, \mathbf{u})$ ; its  $(i, j)$  element is  $\partial f_i(\mathbf{x}, \mathbf{u})/\partial x_j$ . For a real-valued function such as  $C$ , the gradient  $\nabla C(\mathbf{x})$  is a row vector. The costate and the generalized costate (introduced in section 3) are both row vectors, while all other vectors in the paper are column vectors. The  $L^2$  inner product on  $[0, T]$  is denoted  $\langle \cdot, \cdot \rangle$ .

**2. Objective derivative in Case 1.** In Case 1, it is assumed that the control has the form  $\mathbf{u}(t) = \phi_j(\mathbf{x}(t), t)$  for all  $t \in (s_j, s_{j+1})$ ,  $0 \leq j \leq k$ . With this substitution, the dynamics in (1.1) has the form  $\dot{\mathbf{x}}(t) = \mathbf{f}_j(\mathbf{x}(t), t)$  on  $(s_j, s_{j+1})$ , where  $\mathbf{f}_j(\mathbf{x}, t) = \mathbf{f}(\mathbf{x}, \phi_j(\mathbf{x}, t))$ . Note that  $\mathbf{f}_j$  is viewed as a mapping from  $\mathbb{R}^n \times [0, T]$  to  $\mathbb{R}^n$ . The objective is  $C(\mathbf{x}(T))$ , where  $\mathbf{x}$  is the solution to an initial value problem of the form

$$(2.1) \quad \dot{\mathbf{x}}(t) = \mathbf{F}(\mathbf{x}(t), t), \quad \mathbf{F}(\mathbf{x}, t) = \mathbf{f}_j(\mathbf{x}, t) \text{ for } t \in (s_j, s_{j+1}), \quad \mathbf{x}(0) = \mathbf{x}_0,$$

$0 \leq j \leq k$ . In the switch point algorithm, the goal is to minimize the objective value over the choice of the switching points  $s_1, s_2, \dots, s_k$ . This minimization can be done more efficiently if the gradient of the objective with respect to the switching points is known since superlinearly convergent algorithms such as the conjugate gradient method or a quasi-Newton method could be applied. In Theorem 2.4, a formula is derived for the gradient of  $C$  with respect to  $\mathbf{s}$ . The following three preliminary results are used in the analysis.

LEMMA 2.1. If  $\mathbf{x} : [\sigma_1, \sigma_2] \rightarrow \mathbb{R}^n$  is Lipschitz continuous, then so is  $\|\mathbf{x}(\cdot)\|$  and

$$(2.2) \quad \frac{d}{dt}\|\mathbf{x}(t)\| \leq \|\dot{\mathbf{x}}(t)\|$$

for almost every  $t \in [\sigma_1, \sigma_2]$ .

*Proof.* For any  $s, t \in [\sigma_1, \sigma_2]$ , the triangle inequality gives

$$\|\mathbf{x}(s)\| = \|\mathbf{x}(s) - \mathbf{x}(t) + \mathbf{x}(t)\| \leq \|\mathbf{x}(s) - \mathbf{x}(t)\| + \|\mathbf{x}(t)\|.$$

Rearrange this inequality to obtain

$$(2.3) \quad \|\mathbf{x}(s)\| - \|\mathbf{x}(t)\| \leq \|\mathbf{x}(s) - \mathbf{x}(t)\|.$$

Interchanging  $s$  and  $t$  in (2.3) yields  $\|\mathbf{x}(t)\| - \|\mathbf{x}(s)\| \leq \|\mathbf{x}(s) - \mathbf{x}(t)\|$ . Hence, the absolute value of the difference  $\|\mathbf{x}(s)\| - \|\mathbf{x}(t)\|$  is bounded by  $\|\mathbf{x}(t) - \mathbf{x}(s)\|$ . Since  $\mathbf{x}(\cdot)$  is Lipschitz continuous, then so is  $\|\mathbf{x}(\cdot)\|$ . It follows by Rademacher's theorem that both  $\mathbf{x}(\cdot)$  and  $\|\mathbf{x}(\cdot)\|$  are differentiable almost everywhere. Suppose  $t \in (\sigma_1, \sigma_2)$  is a point of differentiability for both  $\mathbf{x}(\cdot)$  and  $\|\mathbf{x}(\cdot)\|$ , and take  $s = t + \Delta t$  in (2.3). Dividing the resulting inequality by  $\Delta t$  and letting  $\Delta t$  tend to zero yields (2.2).  $\square$

The following result can be deduced from Gronwall's inequality.

LEMMA 2.2. If  $w : [0, T] \rightarrow \mathbb{R}$  is absolutely continuous and for some nonnegative scalars  $a$  and  $b$ ,

$$\dot{w}(t) \leq aw(t) + b \quad \text{for almost every } t \in [0, T],$$

then for all  $t \in [0, T]$ , we have

$$(2.4) \quad w(t) \leq e^{at}(w(0) + bt).$$

The following Lipschitz result is deduced from Lemmas 2.1 and 2.2.

COROLLARY 2.3. Suppose that  $\mathbf{x}$  is an absolutely continuous solution of (2.1) and  $\mathbf{y}$  has the same dynamics but a different initial condition:

$$(2.5) \quad \dot{\mathbf{y}}(t) = \mathbf{F}(\mathbf{y}(t), t), \quad \mathbf{y}(0) = \mathbf{y}_0.$$

If for some  $\rho > 0$  and  $L \geq 0$ , independent of  $j$ ,  $\mathbf{f}_j(\boldsymbol{\chi}, t)$ ,  $0 \leq j \leq k$ , is continuous with respect to  $\boldsymbol{\chi}$  and  $t$  and Lipschitz continuous in  $\boldsymbol{\chi}$ , with Lipschitz constant  $L$ , on the tube

$$(2.6) \quad \{(\boldsymbol{\chi}, t) : t \in [s_j, s_{j+1}] \text{ and } \boldsymbol{\chi} \in \mathcal{B}_\rho(\mathbf{x}(t))\},$$

then for any  $\mathbf{y}_0 \in \mathbb{R}^n$  which satisfies  $e^{LT}\|\mathbf{x}_0 - \mathbf{y}_0\| \leq \rho$ , the initial value problem (2.5) has a solution on  $[0, T]$ , and we have

$$(2.7) \quad \|\mathbf{y}(t) - \mathbf{x}(t)\| \leq e^{Lt}\|\mathbf{y}_0 - \mathbf{x}_0\| \quad \text{for all } t \in [0, T].$$

*Proof.* For  $t = 0$  and  $j = 0$ ,  $\mathbf{y}_0$  is in the interior of a face of the tube (2.6). Due to the Lipschitz continuity of  $\mathbf{F}(\cdot, t)$ , a solution to (2.5) exists for near  $t = 0$ . Define  $w(t) = \|\mathbf{y}(t) - \mathbf{x}(t)\|$ , subtract the differential equations (2.1) and (2.5), and take the norm of each side to obtain

$$(2.8) \quad \|\dot{\mathbf{x}}(t) - \dot{\mathbf{y}}(t)\| \leq Lw(t)$$

for any  $t$  where  $\mathbf{y}$  satisfies (2.5) and  $\mathbf{y}(t)$  lies within the tube around  $\mathbf{x}(t)$  where  $\mathbf{F}(\cdot, t)$  satisfies the Lipschitz continuity property. Any solution  $\mathbf{x}$  of (2.1) and  $\mathbf{y}$  of (2.5) is Lipschitz continuous since their derivatives are bounded. Lemma 2.1 and (2.8) imply that  $\dot{w}(t) \leq Lw(t)$  for  $t$  near 0. Hence, for  $t$  near 0, the bound (2.4) of Lemma 2.2 with  $b = 0$  yields

$$w(t) \leq e^{Lt}w(0) = e^{Lt}\|\mathbf{x}_0 - \mathbf{y}_0\| < \rho \quad \text{when } t < T.$$

For  $t < T$ ,  $\mathbf{y}(t)$  continues to lie in the interior of the tube where  $\mathbf{F}(\cdot, t)$  satisfies the Lipschitz property, which leads to (2.7).  $\square$

In order to analyze the objective derivative with respect to the switch points, we need to make a regularity assumption concerning the functions  $\mathbf{f}_j$  in the initial value problem (2.1) to ensure the stability and uniqueness of solutions to (2.1) as the switching points are perturbed around a given  $\mathbf{s}$ .

*State regularity.* Let  $\mathbf{x}$  denote an absolutely continuous solution to (2.1). It is assumed that there exist constants  $\rho > 0$ ,  $s_j^- \in (s_{j-1}, s_j)$ , and  $s_j^+ \in (s_j, s_{j+1})$ ,  $1 \leq j \leq k$ , such that  $\mathbf{f}_j$  is continuously differentiable on the tube

$$\mathcal{T}_j = \{(\boldsymbol{\chi}, t) : t \in [s_j^-, s_j^+] \text{ and } \boldsymbol{\chi} \in \mathcal{B}_\rho(\mathbf{x}(t))\}, \quad 0 \leq j \leq k,$$

where  $s_0^- = 0$  and  $s_{k+1}^+ = T$ . Moreover,  $\mathbf{f}_j(\boldsymbol{\chi}, t)$  is Lipschitz continuously differentiable with respect to  $\boldsymbol{\chi}$  on  $\mathcal{T}_j$ , uniformly in  $j$  and  $t \in [s_j^-, s_j^+]$ .

**THEOREM 2.4.** *Suppose that  $\mathbf{s}$  is a valid choice for the switching points, that the state regularity property holds, and  $C$  is Lipschitz continuously differentiable in a neighborhood of  $\mathbf{x}(T)$ . Then for  $j = 1, 2, \dots, k$ ,  $C(\mathbf{s})$  is differentiable with respect to  $s_j$  and*

$$(2.9) \quad \frac{\partial C}{\partial s_j}(\mathbf{s}) = H_{j-1}(\mathbf{x}(s_j), \mathbf{p}(s_j), s_j) - H_j(\mathbf{x}(s_j), \mathbf{p}(s_j), s_j),$$

where  $H_j(\mathbf{x}, \mathbf{p}, t) = \mathbf{p}\mathbf{f}_j(\mathbf{x}, t)$ , and the row vector  $\mathbf{p} : [0, T] \rightarrow \mathbb{R}^n$  is the solution to the linear differential equation

$$(2.10) \quad \dot{\mathbf{p}}(t) = -\mathbf{p}(t)\nabla_{\mathbf{x}}\mathbf{F}(\mathbf{x}(t), t), \quad t \in [0, T], \quad \mathbf{p}(T) = \nabla C(\mathbf{x}(T)).$$

*Remark 2.1.* The formula (2.9) for the derivative of the objective  $C$  with respect to a switching point generalizes the result derived in [26, Thm. 2] for a bang-bang control. Note that (2.10) differs from the costate equation in (1.2) since the Jacobian of  $\mathbf{f}$  appears in (1.2), while the Jacobian of  $\mathbf{F}$  appears in (2.10). However, when the control  $\mathbf{u}$  enters the Hamiltonian linearly,  $\nabla_{\mathbf{x}}\mathbf{F}$  approaches  $\nabla_{\mathbf{x}}\mathbf{f}$  as the switch points approach the switch points for an optimal solution of the control problem, and the solution  $\mathbf{p}$  of (2.10) approaches the costate of the optimal solution.

*Proof.* Let  $s$  denote the switching point  $s_j$ . To compute the derivative of the cost with respect to  $s$ , we need to compare the cost associated with  $s$  to the cost gotten when  $s$  is changed to  $s + \Delta s$ . Let  $\mathbf{x}$  and  $\mathbf{y}$  denote the solutions of (2.1) associated with the original and the perturbed switching points, respectively. The dynamics associated with these two solutions are identical except for the time interval  $[s, s + \Delta s]$ . With the definitions  $\mathbf{F}_0 = \mathbf{f}_{j-1}$  and  $\mathbf{F}_1 = \mathbf{f}_j$ , the dynamics for  $t \in [s, s + \Delta s]$  and for  $\Delta s$  sufficiently small are

$$(2.11) \quad \dot{\mathbf{x}}(t) = \mathbf{F}_1(\mathbf{x}(t), t) \quad \text{and} \quad \dot{\mathbf{y}}(t) = \mathbf{F}_0(\mathbf{y}(t), t).$$

The objective value associated with the switching point  $s$  is  $C(\mathbf{x}(T))$ , where  $\mathbf{x}$  is the solution of the initial value problem (2.1). The objective value associated with the switching point  $s + \Delta s$  is  $C(\mathbf{y}(T))$  where  $\mathbf{y}$  has the same dynamics as  $\mathbf{x}$  except on the interval  $[s, s + \Delta s]$ . To evaluate the derivative of the objective with respect to  $s$ , we need to form the ratio  $[C(\mathbf{y}(T)) - C(\mathbf{x}(T))]/\Delta s$ , and then let  $\Delta s$  tend to zero.

Since the dynamics for  $\mathbf{x}$  and  $\mathbf{y}$  are identical except for the interval  $[s, s + \Delta s]$ ,  $\mathbf{x}(t) = \mathbf{y}(t)$  for  $t \leq s$ . Let  $\mathbf{x}_s$  denote  $\mathbf{x}(s) = \mathbf{y}(s)$ . Expanding  $\mathbf{x}$  and  $\mathbf{y}$  in Taylor series around  $\mathbf{x}_s$  yields

$$(2.12) \quad \mathbf{x}(s + \Delta s) = \mathbf{x}_s + (\Delta s)\mathbf{F}_1(\mathbf{x}_s, s) + \mathcal{O}((\Delta s)^2),$$

$$(2.13) \quad \mathbf{y}(s + \Delta s) = \mathbf{x}_s + (\Delta s)\mathbf{F}_0(\mathbf{x}_s, s) + \mathcal{O}((\Delta s)^2),$$

where the remainder term is  $\mathcal{O}((\Delta s)^2)$  due to the state regularity property, which ensures the Lipschitz continuous differentiability of  $\mathbf{F}_0$  and  $\mathbf{F}_1$ . For  $t > s + \Delta s$ ,  $\dot{\mathbf{x}}(t) = \mathbf{F}(\mathbf{x}(t), t)$  and  $\dot{\mathbf{y}}(t) = \mathbf{F}(\mathbf{y}(t), t)$  since the dynamics of  $\mathbf{x}$  and  $\mathbf{y}$  only differ on  $[s, s + \Delta s]$ . Apply Corollary 2.3 to the interval  $[s + \Delta s, T]$ . Based on the expressions (2.12) and (2.13) for the values of  $\mathbf{x}$  and  $\mathbf{y}$  at  $s + \Delta s$ , it follows from Corollary 2.3 that

$$(2.14) \quad \|\mathbf{y}(t) - \mathbf{x}(t)\| = \mathcal{O}(\Delta s) \quad \text{for all } t \in [s + \Delta s, T].$$

Let  $\mathbf{z} : [s + \Delta s, T] \rightarrow \mathbb{R}^n$  be the solution to the linear differential equation

$$(2.15) \quad \dot{\mathbf{z}}(t) = \nabla_x \mathbf{F}(\mathbf{x}(t), t)\mathbf{z}(t), \quad \mathbf{z}(s + \Delta s) = \Delta s[\mathbf{F}_0(\mathbf{x}_s, s) - \mathbf{F}_1(\mathbf{x}_s, s)].$$

By assumption,  $\mathbf{x}$  is absolutely continuous and hence bounded; consequently, there exists a scalar  $a$  such that  $\|\nabla_x \mathbf{F}(\mathbf{x}, t)\| \leq a$  for all  $t \in [s, T]$  and  $\mathbf{x} \in \mathcal{B}_\rho(\mathbf{x}(t))$ . Take the norm of the differential equation for  $\mathbf{z}$  and apply Lemma 2.1 to obtain

$$\frac{d\|\mathbf{z}(t)\|}{dt} \leq a\|\mathbf{z}(t)\| \quad \text{for all } t \in [s + \Delta s, T].$$

By Lemma 2.2 with  $b = 0$  and the specified initial condition in (2.15), it follows that

$$(2.16) \quad \|\mathbf{z}(t)\| = \mathcal{O}(\Delta s) \quad \text{for all } t \in [s + \Delta s, T].$$

Define  $\delta(t) = \mathbf{y}(t) - \mathbf{x}(t) - \mathbf{z}(t)$  for every  $t \in [s + \Delta s, T]$  and

$$\mathbf{x}(\alpha, t) = \mathbf{x}(t) + \alpha(\mathbf{y}(t) - \mathbf{x}(t)).$$

Differentiating  $\delta$  and utilizing a Taylor expansion with integral remainder term, we obtain for all  $t \in [s + \Delta s, T]$ ,

$$(2.17) \quad \begin{aligned} \dot{\delta}(t) &= \dot{\mathbf{y}}(t) - \dot{\mathbf{x}}(t) - \dot{\mathbf{z}}(t) = \mathbf{F}(\mathbf{y}(t), t) - \mathbf{F}(\mathbf{x}(t), t) - \nabla_x \mathbf{F}(\mathbf{x}(t), t)\mathbf{z}(t) \\ &= \left( \int_0^1 \nabla_x \mathbf{F}(\mathbf{x}(\alpha, t), t) d\alpha \right) (\mathbf{y}(t) - \mathbf{x}(t)) - \left( \int_0^1 \nabla_x \mathbf{F}(\mathbf{x}(t), t) d\alpha \right) \mathbf{z}(t) \\ &= \left( \int_0^1 [\nabla_x \mathbf{F}(\mathbf{x}(\alpha, t), t) - \nabla_x \mathbf{F}(\mathbf{x}(t), t)] d\alpha \right) \mathbf{z}(t) + \left( \int_0^1 \nabla_x \mathbf{F}(\mathbf{x}(\alpha, t), t) d\alpha \right) \delta(t). \end{aligned}$$

Take  $\Delta s$  in (2.14) small enough that  $\mathbf{y}(t)$  lies in the tube around  $\mathbf{x}(t)$  where  $\nabla_x \mathbf{F}$  is Lipschitz continuous. By the definition of  $a$ , we have  $\|\nabla_x \mathbf{F}(\mathbf{x}(\alpha, t), t)\| \leq a$ . If  $L$  is the Lipschitz constant for  $\nabla_x \mathbf{F}$ , then by (2.14), we have

$$\|\nabla_x \mathbf{F}(\mathbf{x}(\alpha, t), t) - \nabla_x \mathbf{F}(\mathbf{x}(t), t)\| \leq \alpha L \|\mathbf{y}(t) - \mathbf{x}(t)\| = \mathcal{O}(\Delta s).$$



Hence, taking the norm of each side of (2.17) and utilizing Lemma 2.1 on the left side and the triangle inequality and the bound (2.16) on the right side yields

$$(2.18) \quad \frac{d\|\delta(t)\|}{dt} \leq a\|\delta(t)\| + \mathcal{O}((\Delta s)^2).$$

By (2.12), (2.13), and (2.15), we have

$$\delta(s + \Delta s) = \mathbf{y}(s + \Delta s) - \mathbf{x}(s + \Delta s) - \mathbf{z}(s + \Delta s) = \mathcal{O}((\Delta s)^2).$$

Consequently, by (2.18) and Lemma 2.2, we deduce that

$$(2.19) \quad \|\delta(t)\| = \mathcal{O}((\Delta s)^2) \quad \text{for all } t \in [s + \Delta s, T].$$

If  $\mathbf{p}$  is the solution of (2.10) and  $\mathbf{z}$  is the solution of (2.15), then we have

$$\begin{aligned} 0 &= \int_{s+\Delta s}^T \mathbf{p}(t) \left[ \nabla_x \mathbf{F}(\mathbf{x}(t), t) \mathbf{z}(t) - \dot{\mathbf{z}}(t) \right] dt \\ &= \int_{s+\Delta s}^T \left[ \mathbf{p}(t) \nabla_x \mathbf{F}(\mathbf{x}(t), t) + \dot{\mathbf{p}}(t) \right] \mathbf{z}(t) dt - \mathbf{p}(T) \mathbf{z}(T) + \mathbf{p}(s + \Delta s) \mathbf{z}(s + \Delta s) \\ &= -\mathbf{p}(T) \mathbf{z}(T) + \mathbf{p}(s + \Delta s) \mathbf{z}(s + \Delta s) \\ (2.20) \quad &= \Delta s \mathbf{p}(s + \Delta s) (\mathbf{F}_0(\mathbf{x}_s, s) - \mathbf{F}_1(\mathbf{x}_s, s)) - \nabla C(\mathbf{x}(T)) \mathbf{z}(T). \end{aligned}$$

Since  $C$  is Lipschitz continuously differentiable in a neighborhood of  $\mathbf{x}(T)$ , the difference between the perturbed objective and the original objective can be expressed

$$(2.21) \quad C(\mathbf{y}(T)) - C(\mathbf{x}(T)) = \nabla C(\mathbf{x}_\Delta)(\mathbf{y}(T) - \mathbf{x}(T)),$$

where  $\mathbf{x}_\Delta \in [\mathbf{y}(T), \mathbf{x}(T)]$ ; that is,  $\mathbf{x}_\Delta$  is a point on the line segment connecting  $\mathbf{y}(T)$  and  $\mathbf{x}(T)$ . Since the distance between  $\mathbf{x}(t)$  and  $\mathbf{y}(t)$  is  $\mathcal{O}(\Delta s)$  by (2.14), the distance between  $\mathbf{x}_\Delta$  and  $\mathbf{x}(T)$  is  $\mathcal{O}(\Delta s)$ . Add the right side of (2.20) to the right side of (2.21) and substitute

$$\mathbf{y}(T) - \mathbf{x}(T) = \mathbf{y}(T) - \mathbf{x}(T) - \mathbf{z}(T) + \mathbf{z}(T) = \delta(T) + \mathbf{z}(T)$$

to obtain

$$(2.22) \quad \begin{aligned} C(\mathbf{y}(T)) - C(\mathbf{x}(T)) &= \nabla C(\mathbf{x}_\Delta) \delta(T) + [\nabla C(\mathbf{x}_\Delta) - \nabla C(\mathbf{x}(T))] \mathbf{z}(T) \\ &\quad + \Delta s \mathbf{p}(s + \Delta s) [\mathbf{F}_0(\mathbf{x}_s, s) - \mathbf{F}_1(\mathbf{x}_s, s)]. \end{aligned}$$

By (2.19),  $\|\delta(T)\| = \mathcal{O}((\Delta s)^2)$  so  $|\nabla C(\mathbf{x}_\Delta) \delta(T)| = \mathcal{O}((\Delta s)^2)$ . Since  $C$  is Lipschitz continuously differentiable at  $\mathbf{x}(T)$ , the distance from  $\mathbf{x}_\Delta$  to  $\mathbf{x}(T)$  is at most  $\mathcal{O}(\Delta s)$  by (2.14), and  $\mathbf{z}(T) = \mathcal{O}(\Delta s)$  by (2.16), we have

$$\|[\nabla C(\mathbf{x}_\Delta) - \nabla C(\mathbf{x}(T))] \mathbf{z}(T)\| = \mathcal{O}((\Delta s)^2).$$

Consequently, the first two terms on the right side of (2.22) are  $\mathcal{O}((\Delta s)^2)$ . Divide (2.22) by  $\Delta s$  and let  $\Delta s$  tend to zero to obtain

$$\lim_{\Delta s \rightarrow 0} \frac{C(\mathbf{y}(T)) - C(\mathbf{x}(T))}{\Delta s} = \mathbf{p}(s) [\mathbf{F}_0(\mathbf{x}_s, s) - \mathbf{F}_1(\mathbf{x}_s, s)],$$

which completes the proof since  $\mathbf{F}_0 = \mathbf{f}_{j-1}$  and  $\mathbf{F}_1 = \mathbf{f}_j$ .  $\square$

**3. Objective derivative in Case 2.** In Case 2, an optimal control which minimizes the Hamiltonian also depends on  $\mathbf{p}$  in the singular region, so the control has the form  $\mathbf{u}(t) = \phi_j(\mathbf{x}(t), \mathbf{p}(t), t)$  for  $t \in (s_j, s_{j+1})$ . The functions  $\mathbf{f}_j$  and  $\mathbf{F}$  of section 2 now have the form  $\mathbf{f}_j(\mathbf{x}, \mathbf{p}, t) = \mathbf{f}(\mathbf{x}, \phi_j(\mathbf{x}, \mathbf{p}, t))$  and  $\mathbf{F}(\mathbf{x}, \mathbf{p}, t) = \mathbf{f}_j(\mathbf{x}, \mathbf{p}, t)$  for  $t \in (s_j, s_{j+1})$ . The Jacobian appearing in the costate dynamics is  $\nabla_x \mathbf{f}(\mathbf{x}, \mathbf{u})$ ; in the switch point algorithm, we evaluate this dynamics at  $\mathbf{u} = \phi_j(\mathbf{x}, \mathbf{p}, t)$ . Hence, analogous to the definition given in section 2, let us define

$$\mathbf{f}_{jx}(\mathbf{x}, \mathbf{p}, t) = \nabla_x \mathbf{f}(\mathbf{x}, \mathbf{u}) \Big|_{\mathbf{u} = \phi_j(\mathbf{x}, \mathbf{p}, t)}.$$

Also, we define

$$\mathbf{F}_x(\mathbf{x}, \mathbf{p}, t) = \mathbf{f}_{jx}(\mathbf{x}, \mathbf{p}, t) \quad \text{for } t \in (s_j, s_{j+1}).$$

In the switch point algorithm, we consider the coupled system

$$(3.1) \quad \dot{\mathbf{x}}(t) = \mathbf{F}(\mathbf{x}(t), \mathbf{p}(t), t), \quad \dot{\mathbf{p}}(t) = -\mathbf{p}(t)\mathbf{F}_x(\mathbf{x}(t), \mathbf{p}(t), t),$$

where  $(\mathbf{x}(0), \mathbf{p}(0)) = (\mathbf{x}_0, \mathbf{p}_0)$ . If  $\mathbf{u}^*$  is a local minimizer for the control problem (1.1) and  $(\mathbf{x}^*, \mathbf{p}^*)$  are the associated state and costate, then we could recover  $\mathbf{u}^*(t) = \phi_j(\mathbf{x}^*(t), \mathbf{p}^*(t), t)$ ,  $t \in (s_j, s_{j+1})$ , by integrating the coupled system (3.1) forward in time starting from the initial condition  $(\mathbf{x}(0), \mathbf{p}(0)) = (\mathbf{x}_0, \mathbf{p}^*(0))$ . From this perspective, we can think of the objective  $C(\mathbf{x}(T))$  as being a function  $C(\mathbf{s}, \mathbf{p}_0)$  that depends on both the switching points and the starting value  $\mathbf{p}_0$  for the costate (the starting condition for the state  $\mathbf{x}_0$  is given). To solve the control problem, we will search for a local minimizer of  $C(\mathbf{s}, \mathbf{p}_0)$ . Again, to exploit superlinearly convergent optimization algorithms, the derivatives of  $C$  with respect to both  $\mathbf{s}$  and  $\mathbf{p}_0$  should be evaluated.

The derivative of the objective with respect to the switching points in Case 2 is a corollary of Theorem 2.4. In this case, the  $\mathbf{x}$  that satisfies (2.1) is identified with the pair  $(\mathbf{x}, \mathbf{p})$  which solves the coupled system (3.1). The pair  $(\mathbf{x}, \mathbf{p})$  might be viewed as a *generalized state* in the sense that for a given starting value  $\mathbf{p}(0) = \mathbf{p}_0$  and for a given choice of the switch points, we can in principle integrate forward in time the coupled system (3.1) to evaluate the objective  $C(\mathbf{x}(T))$ . The generalized version of the state regularity property, which applies to the pair  $(\mathbf{x}, \mathbf{p})$ , is the following.

*Generalized state regularity.* Let  $(\mathbf{x}, \mathbf{p})$  denote an absolutely continuous solution to (3.1). It is assumed that there exist constants  $\rho > 0$ ,  $s_j^- \in (s_{j-1}, s_j)$ , and  $s_j^+ \in (s_j, s_{j+1})$ ,  $1 \leq j \leq k$ , such that the pair  $(\mathbf{f}_j, \mathbf{f}_{jx})$  is continuously differentiable on the tube

$$\mathcal{T}_j = \{(\boldsymbol{\chi}, t) : t \in [s_j^-, s_{j+1}^+] \text{ and } \boldsymbol{\chi} \in \mathcal{B}_\rho(\mathbf{x}(t))\}, \quad 0 \leq j \leq k,$$

where  $s_0^- = 0$  and  $s_{k+1}^+ = T$ . Moreover,  $(\mathbf{f}_j(\boldsymbol{\chi}, t), \mathbf{f}_{jx}(\boldsymbol{\chi}, t))$  is Lipschitz continuously differentiable with respect to  $\boldsymbol{\chi}$  on  $\mathcal{T}_j$ , uniformly in  $j$  and  $t \in [s_j^-, s_j^+]$ .

The *generalized costate* associated with the system (3.1) is a row vector  $\mathbf{y} \in \mathbb{R}^{2n}$  whose first  $n$  components are denoted  $\mathbf{y}_1$  and whose second  $n$  components are denoted  $\mathbf{y}_2$ . The generalized Hamiltonian is defined by

$$\mathcal{H}_j(\mathbf{x}, \mathbf{p}, \mathbf{y}, t) = \mathbf{y}_1 \mathbf{f}_j(\mathbf{x}, \mathbf{p}, t) - \mathbf{p} \mathbf{f}_{jx}(\mathbf{x}, \mathbf{p}, t) \mathbf{y}_2^\top, \quad 0 \leq j \leq k.$$

The generalized costate  $\mathbf{y} : [0, T] \rightarrow \mathbb{R}^{2n}$  is the solution of the linear system of differential equations

$$(3.2) \quad \dot{\mathbf{y}}_1(t) = -\nabla_x \mathcal{H}_j(\mathbf{x}(t), \mathbf{p}(t), \mathbf{y}(t), t), \quad \mathbf{y}_1(T) = \nabla C(\mathbf{x}(T)),$$

$$(3.3) \quad \dot{\mathbf{y}}_2(t) = -\nabla_p \mathcal{H}_j(\mathbf{x}(t), \mathbf{p}(t), \mathbf{y}(t), t), \quad \mathbf{y}_2(T) = \mathbf{0},$$

on  $(s_j, s_{j+1})$  for  $j = k, k-1, \dots, 0$ . If the generalized state regularity property holds, then by Theorem 2.4, we have

$$(3.4) \quad \frac{\partial C}{\partial s_j}(\mathbf{s}, \mathbf{p}_0) = \mathcal{H}_{j-1}(\mathbf{x}(s_j), \mathbf{p}(s_j), \mathbf{y}(s_j), s_j) - \mathcal{H}_j(\mathbf{x}(s_j), \mathbf{p}(s_j), \mathbf{y}(s_j), s_j)$$

for  $j = 1, 2, \dots, k$ . Note that the boundary condition for  $\mathbf{y}_2$  is  $\mathbf{y}_2(T) = \mathbf{0}$  since there is no  $\mathbf{p}$  in the objective, and the objective  $C$  only depends on  $\mathbf{x}(T)$ .

*Remark 3.1.* If the formula (3.4) is used to evaluate the derivative of the objective with respect to a switch point in the case where  $\phi$  does not depend on  $\mathbf{p}$ , then the formula (3.4) reduces to the formula (2.9). In particular, when  $\phi$  does not depend on  $\mathbf{p}$ , the dynamics for  $\mathbf{y}_2$  becomes

$$\dot{\mathbf{y}}_2(t) = \mathbf{y}_2 \mathbf{F}_x(\mathbf{x}, \phi_j(\mathbf{x}(t), t))^T, \quad t \in (s_j, s_{j+1}), \quad j = k, k-1, \dots, 0, \quad \mathbf{y}_2(T) = \mathbf{0}.$$

Consequently,  $\mathbf{y}_2$  is the solution to a linear differential equation with the initial condition  $\mathbf{y}_2(T) = \mathbf{0}$ . The unique solution is  $\mathbf{y}_2 = \mathbf{0}$ , and when  $\mathbf{y}_2 = \mathbf{0}$ , (3.2) is the same as (2.10). Thus  $\mathbf{y}_1 = \mathbf{p}$  and the formula (3.4) is the same as (2.9).

Now let us consider the gradient of  $C(\mathbf{s}, \mathbf{p}_0)$  with respect to  $\mathbf{p}_0$ . Let  $(\mathbf{x}, \mathbf{p})$  denote a solution of (3.1) for a given starting condition  $\mathbf{p}(0) = \mathbf{p}_0$ , and let  $(\bar{\mathbf{x}}, \bar{\mathbf{p}})$  denote a solution corresponding to  $\mathbf{p}(0) = \bar{\mathbf{p}}_0$ . Let  $\mathbf{y}$  denote the generalized costate associated with  $(\bar{\mathbf{x}}, \bar{\mathbf{p}})$ . Since  $(\mathbf{x}, \mathbf{p})$  is a solution of (3.1), we have

$$(3.5) \quad 0 = \int_0^T \mathbf{y}_1(t) [\mathbf{F}(\mathbf{x}(t), \mathbf{p}(t), t) - \dot{\mathbf{x}}(t)] - [\mathbf{p}(t) \mathbf{F}_x(\mathbf{x}(t), \mathbf{p}(t), t) + \dot{\mathbf{p}}(t)] \mathbf{y}_2^T(t) dt.$$

The two derivative terms in (3.5) are integrated by parts to obtain

$$(3.6) \quad -[\langle \mathbf{y}_1, \dot{\mathbf{x}} \rangle + \langle \mathbf{y}_2, \dot{\mathbf{p}} \rangle] = \langle \dot{\mathbf{y}}_1, \mathbf{x} \rangle + \langle \dot{\mathbf{y}}_2, \mathbf{p} \rangle + \mathbf{y}_1(0) \mathbf{x}_0 + \mathbf{y}_2(0) \mathbf{p}_0^T - \nabla C(\bar{\mathbf{x}}(T)) \mathbf{x}(T),$$

where  $\langle \cdot, \cdot \rangle$  denotes the  $L^2$  inner product on  $[0, T]$ , and the boundary conditions in (3.1), (3.2), and (3.3) are used to simplify the boundary terms.

We now combine (3.5) and (3.6), differentiate the resulting identity with respect to  $\mathbf{p}_0$ , and evaluate the derivative at  $\mathbf{p}_0 = \bar{\mathbf{p}}_0$ . Recall that  $\mathbf{y}$  is independent of  $\mathbf{p}_0$  since it corresponds to (3.2) and (3.3) in the special case where  $(\mathbf{x}, \mathbf{p}) = (\bar{\mathbf{x}}, \bar{\mathbf{p}})$ . The only terms depending on  $\mathbf{p}_0$  are those involving  $\mathbf{x}$  and  $\mathbf{p}$ , the solution of (3.1). In particular, the partial derivatives of the three boundary terms  $\mathbf{y}_1(0) \mathbf{x}_0$ ,  $\mathbf{y}_2(0) \mathbf{p}_0^T$ , and  $\nabla C(\bar{\mathbf{x}}(T)) \mathbf{x}(T)$  with respect to  $\mathbf{p}_0$  are zero,  $\mathbf{y}_2(0)$ , and  $\nabla C(\bar{\mathbf{x}}(T)) \partial \mathbf{x}(T) / \partial \mathbf{p}_0$ , respectively. When we differentiate (3.5) and (3.6) with respect to  $\mathbf{p}_0$  and evaluate at  $\mathbf{p}_0 = \bar{\mathbf{p}}_0$ , every term cancels except for these three terms (and one of these three terms is zero). To see how these terms in the integrals cancel, let us consider those terms with the common factor  $(\partial \mathbf{x} / \partial \mathbf{p}_0)(t)$ . This factor in the integral is multiplied by

$$\dot{\mathbf{y}}_1(t) + \nabla_x \mathcal{H}_j(\mathbf{x}(t), \mathbf{p}(t), \mathbf{y}(t), t),$$

which vanishes for  $\mathbf{x} = \bar{\mathbf{x}}$  and  $\mathbf{p} = \bar{\mathbf{p}}$  by (3.2). The remaining terms with the common factor  $(\partial \mathbf{p} / \partial \mathbf{p}_0)(t)$  vanish due (3.3). After taking into account the three boundary terms in (3.6), we obtain

$$(3.7) \quad \left. \frac{\partial C(\mathbf{s}, \mathbf{p}_0)}{\partial \mathbf{p}_0} \right|_{\mathbf{p}_0 = \bar{\mathbf{p}}_0} = \nabla C(\bar{\mathbf{x}}(T)) \left. \frac{\partial \mathbf{x}(T)}{\partial \mathbf{p}_0} \right|_{\mathbf{p}_0 = \bar{\mathbf{p}}_0} = \mathbf{y}_2(0).$$

In summary, the gradient of the objective with respect to  $\mathbf{p}_0$  is available almost for free after evaluating the derivative of the objective with respect to the switching points; the gradient is simply  $\mathbf{y}_2(0)$ . As pointed out in Remark 3.1,  $\mathbf{y}_2 = \mathbf{0}$  when the  $\phi_j$  are independent of  $\mathbf{p}$ .

**4. Free terminal time.** So far, the terminal time  $T$  has been fixed. Let us now suppose that the terminal time is free, and we are minimizing over both the terminal time  $T$  and over the control  $\mathbf{u}$ . It is assumed that the control constraint set  $\mathcal{U}$  is independent of  $t$ , and we make the change of variable  $t = \tau T$ , where  $0 \leq \tau \leq 1$ . After the change of variables, both the state and the control are functions of  $\tau$  rather than  $t$ . The reformulated optimization problem is

$$(4.1) \quad \min C(\mathbf{x}(1)) \quad \text{subject to} \quad \dot{\mathbf{x}}(\tau) = T\mathbf{f}(\mathbf{x}(\tau), \mathbf{u}(\tau)), \quad \mathbf{x}(0) = \mathbf{x}_0, \quad \mathbf{u}(\tau) \in \mathcal{U},$$

where  $\mathbf{x} : [0, 1] \rightarrow \mathbb{R}^n$  is absolutely continuous and  $\mathbf{u} : [0, 1] \rightarrow \mathbb{R}^m$  is essentially bounded. In the reformulated problem, the free terminal time  $T$  appears as a parameter in the system dynamics.

For fixed  $T$ , the optimization problem over the control has the same structure as that of the problem analyzed in sections 2 and 3. Hence, the previously derived formula for the derivative of the objective with respect to a switching point remains applicable. If a gradient-based algorithm will be used to solve (4.1), then we also need a formula for the derivative of the objective with respect to  $T$  when the switch points for the control are fixed. Since the switch points for the control are fixed throughout this section, the objective value in (4.1) only depends on the choice of the parameter  $T$  in the dynamics. Assuming that for some given  $T$  there exists a solution  $\mathbf{x}$  to the dynamics in (4.1), we let  $C(T) := C(\mathbf{x}(1))$  denote the objective value. By the chain rule,

$$(4.2) \quad \frac{dC(T)}{dT} = \nabla C[\mathbf{x}(1)] \frac{d\mathbf{x}(1)}{dT}.$$

Similar to the approach in section 3, our goal is to obtain an expression for the right side of (4.2) that avoids the computation of the derivative of the state with respect to  $T$ . Let us first consider Case 1 where the control has the form  $\mathbf{u}(\tau) = \phi_j(\mathbf{x}(\tau), \tau)$  for all  $\tau \in (s_j, s_{j+1})$ ,  $0 \leq j \leq k$ .

**THEOREM 4.1.** *Suppose that for  $T = \bar{T}$ ,  $\mathbf{x} = \bar{\mathbf{x}}$  is an absolutely continuous solution of*

$$(4.3) \quad \dot{\mathbf{x}}(\tau) = T\mathbf{F}(\mathbf{x}(\tau), \tau), \quad \mathbf{x}(0) = \mathbf{x}_0, \quad 0 \leq \tau \leq 1,$$

where  $\mathbf{F}$  is defined in (2.1). We assume that for some  $\rho > 0$ ,  $\mathbf{f}_j(\boldsymbol{\chi}, \tau)$ ,  $0 \leq j \leq k$ , is continuous with respect to  $\boldsymbol{\chi}$  and  $\tau$  and Lipschitz continuous with respect to  $\boldsymbol{\chi}$  on the tube

$$\{(\boldsymbol{\chi}, \tau) : \tau \in [s_j, s_{j+1}] \text{ and } \boldsymbol{\chi} \in \mathcal{B}_\rho(\mathbf{x}(\tau))\}, \quad 0 \leq j \leq k.$$

Then we have

$$(4.4) \quad \left. \frac{dC(T)}{dT} \right|_{T=\bar{T}} = \int_0^1 H(\bar{\mathbf{x}}(\tau), \mathbf{p}(\tau), \tau) d\tau,$$

where  $H(\mathbf{x}, \mathbf{p}, \tau) = H_j(\mathbf{x}, \mathbf{p}, \tau)$  when  $s_j \leq \tau \leq s_{j+1}$ , and the row vector  $\mathbf{p} : [0, 1] \rightarrow \mathbb{R}^n$  is the solution to the linear differential equation

$$(4.5) \quad \dot{\mathbf{p}}(\tau) = -\bar{T}\mathbf{p}(\tau)\nabla_{\mathbf{x}}\mathbf{F}(\bar{\mathbf{x}}(\tau), \tau), \quad \tau \in [0, 1], \quad \mathbf{p}(1) = \nabla C[\bar{\mathbf{x}}(1)].$$

*Proof.* If  $\mathbf{p}$  denotes the costate given by (4.5), and  $\mathbf{x}$  for  $T$  near  $\bar{T}$  denotes the solution of (4.3), then we have the identity

$$(4.6) \quad \begin{aligned} 0 &= \int_0^1 \mathbf{p}(\tau) [T\mathbf{F}(\mathbf{x}(\tau), \tau) - \dot{\mathbf{x}}(\tau)] d\tau \\ &= \mathbf{p}(0)\mathbf{x}_0 - \mathbf{p}(1)\mathbf{x}(1) + \int_0^1 [T\mathbf{p}(\tau)\mathbf{F}(\mathbf{x}(\tau), \tau) + \dot{\mathbf{p}}(\tau)\mathbf{x}(\tau)] d\tau, \end{aligned}$$

where the second equation comes from an integration by parts. Let us differentiate with respect to  $T$ . Since  $\mathbf{p}$  in (4.5) does not depend on  $T$ , the derivative of  $\mathbf{p}(0)\mathbf{x}_0$  with respect to  $T$  is zero. Hence, the derivative of (4.6) with respect to  $T$ , evaluated at  $T = \bar{T}$ , yields

$$\mathbf{p}(1) \frac{d\mathbf{x}(1)}{dT} \Big|_{T=\bar{T}} = \int_0^1 \left\{ \mathbf{p}(\tau)\mathbf{F}(\bar{\mathbf{x}}(\tau), \tau) + [\bar{T}\nabla_{\mathbf{x}}\mathbf{F}(\bar{\mathbf{x}}(\tau), \tau) + \dot{\mathbf{p}}(\tau)] \frac{d\mathbf{x}(\tau)}{dT} \Big|_{T=\bar{T}} \right\} d\tau.$$

Substituting for  $\mathbf{p}$  from (4.5), the factor multiplying  $d\mathbf{x}(\tau)/dT$  is zero. It follows from (4.2) that

$$\frac{dC(T)}{dT} \Big|_{T=\bar{T}} = \nabla C[\bar{\mathbf{x}}(1)] \frac{d\mathbf{x}(1)}{dT} \Big|_{T=\bar{T}} = \int_0^1 H(\bar{\mathbf{x}}(\tau), \mathbf{p}(\tau), \tau) d\tau,$$

which completes the proof.  $\square$

*Remark 4.1.* Note that the Hamiltonian in (4.4) does not contain the terminal time, while the Hamiltonian associated with (4.3) and the objective derivative with respect to a switch point does contain the terminal time. It follows from (4.4) that the integral of the Hamiltonian vanishes along an optimal solution to the control problem (1.1). Hence, due to the constancy of the Hamiltonian along an optimal solution, (4.4) implies that the Hamiltonian vanishes along an optimal solution, a classic first-order optimality condition for free terminal time control problems.

Now consider Case 2 where the control has the form  $\mathbf{u}(\tau) = \phi_j(\mathbf{x}(\tau), \mathbf{p}(\tau), \tau)$  for all  $\tau \in (s_j, s_{j+1})$ ,  $0 \leq j \leq k$ . The generalized state  $(\mathbf{x}, \mathbf{p})$  satisfies the coupled system

$$(4.7) \quad \dot{\mathbf{x}}(t) = T\mathbf{F}(\mathbf{x}(t), \mathbf{p}(t), t), \quad \dot{\mathbf{p}}(t) = -T\mathbf{p}(t)\mathbf{F}_x(\mathbf{x}(t), \mathbf{p}(t), t),$$

where  $(\mathbf{x}(0), \mathbf{p}(0)) = (\mathbf{x}_0, \mathbf{p}_0)$  and the terminal time  $T$  is a parameter in the equations. For fixed  $T$  and a fixed choice of the switching times  $\mathbf{s}$ , the derivative of the objective with respect to  $\mathbf{p}_0$  is given by the formula (3.7). Now for a fixed choice of both  $\mathbf{s}$  and  $\mathbf{p}_0$ , say,  $\mathbf{p}_0 = \bar{\mathbf{p}}_0$ , our goal is to evaluate the objective derivative with respect to  $T$  evaluated at some given terminal time  $T = \bar{T}$ , assuming a solution  $(\bar{\mathbf{x}}, \bar{\mathbf{p}})$  to (4.7) exists for  $T = \bar{T}$  and  $\mathbf{p}_0 = \bar{\mathbf{p}}_0$ .

The analogue of (4.6) is a slightly modified version of (3.5) where the integration limit  $T$  is replaced by 1, while  $T$  appears as a parameter next to the dynamics:

$$(4.8) \quad 0 = \int_0^1 \mathbf{y}_1(t) [T\mathbf{F}(\mathbf{x}(t), \mathbf{p}(t), t) - \dot{\mathbf{x}}(t)] - [T\mathbf{p}(t)\mathbf{F}_x(\mathbf{x}(t), \mathbf{p}(t), t) + \dot{\mathbf{p}}(t)] \mathbf{y}_2^T(t) dt,$$

where  $(\mathbf{x}, \mathbf{p})$  is the solution to (4.7) corresponding to a general  $T$ , but with the initial condition  $\mathbf{p}_0 = \bar{\mathbf{p}}_0$ . The generalized costate  $\mathbf{y}$  in (4.8) is the solution to

$$\begin{aligned} \dot{\mathbf{y}}_1(t) &= -\bar{T}\nabla_{\mathbf{x}}\mathcal{H}_j(\bar{\mathbf{x}}(t), \bar{\mathbf{p}}(t), \mathbf{y}(t), t), & \mathbf{y}_1(T) &= \nabla C(\bar{\mathbf{x}}(T)), \\ \dot{\mathbf{y}}_2(t) &= -\bar{T}\nabla_{\mathbf{p}}\mathcal{H}_j(\bar{\mathbf{x}}(t), \bar{\mathbf{p}}(t), \mathbf{y}(t), t), & \mathbf{y}_2(T) &= \mathbf{0}, \end{aligned}$$

on  $(s_j, s_{j+1})$  for  $j = k, k-1, \dots, 0$ . Proceeding as in the proof of Theorem 4.1, we integrate by parts in (4.8), differentiate with respect to  $T$ , and evaluate at  $T = \bar{T}$  to obtain the relation

$$\left. \frac{dC(T)}{dT} \right|_{T=\bar{T}} = \nabla C[\bar{\mathbf{x}}(1)] \left. \frac{d\mathbf{x}(1)}{dT} \right|_{T=\bar{T}} = \int_0^1 \mathcal{H}(\bar{\mathbf{x}}(\tau), \bar{\mathbf{p}}(\tau), \mathbf{y}(\tau), \tau) d\tau,$$

where  $\mathcal{H}(\mathbf{x}, \mathbf{p}, \mathbf{y}, \tau) = \mathcal{H}_j(\mathbf{x}, \mathbf{p}, \mathbf{y}, \tau)$  when  $s_j \leq \tau \leq s_{j+1}$ ,  $0 \leq j \leq k$ .

**5. Starting guess.** This section discusses how to generate a starting guess for the switch point algorithm. Detailed illustrations of these techniques are given in [6]. If the optimal control is bang-bang without singular intervals, then the optimization problem could be discretized by Euler's method, and the location of the switching point can often be estimated with a few iterations of a gradient or a conjugate gradient method. On the other hand, when a singular interval is present, wild oscillations in the control can occur and the problem becomes more difficult. An effective way to approximate the optimal control in the singular setting is to incorporate TV regularization in the objective. TV regularization has been very effective in image restoration since it preserves sharp edges; for singular optimal control problems, it helps to remove the wild oscillations in the control, and better exposes the switch points.

We consider the Euler discretization of (1.1) with  $\rho$ -amplified TV regularization:

$$(5.1) \quad \min C(\mathbf{x}_N) + \rho \sum_{i=1}^m \sum_{j=1}^{N-1} |u_{ij} - u_{i,j-1}|$$

subject to  $\mathbf{x}_{j+1} = \mathbf{x}_j + h\mathbf{f}(\mathbf{x}_j, \mathbf{u}_j)$ ,  $\mathbf{u}_j \in \mathcal{U}(t_j)$ ,

where  $0 \leq j \leq N-1$ ,  $h = T/N$ ,  $t_j = jh$ , and  $N$  is the number of mesh intervals. The parameter  $\rho$  controls the strength of the TV regularization term, and as  $\rho$  increases, the oscillations in  $u$  should decrease. The nonsmooth problem (5.1) is equivalent to the smooth optimization problem

$$(5.2) \quad \min C(\mathbf{x}_N) + \rho \sum_{i=1}^m \sum_{j=1}^{N-2} v_{ij} + w_{ij}$$

s.t.  $\mathbf{x}_{j+1} = \mathbf{x}_j + h\mathbf{f}(\mathbf{x}_j, \mathbf{u}_j)$ ,  $\mathbf{u}_j \in \mathcal{U}(t_j)$ ,  $\mathbf{u}_{l+1} - \mathbf{u}_l = \mathbf{v}_l - \mathbf{w}_l$ ,  $\mathbf{v}_l \geq \mathbf{0}$ ,  $\mathbf{w}_l \geq \mathbf{0}$ ,

where  $0 \leq j \leq N-1$  and  $0 \leq l \leq N-2$ . The equivalence between (5.1) and (5.2) is due to the well-known property in optimization that

$$|u| = \min\{v + w : u = v - w, v \geq 0, w \geq 0\}.$$

The smooth TV-regularized problem (5.2) can be solved quickly by the polyhedral active set algorithm (PASA) [27] due to the sparsity of the linear constraints.

Figure 5.1 shows how the optimal control of (5.1) for the catalyst mixing problem of the next section depends on  $\rho$ . When  $\rho = 0$  the control oscillates wildly, when  $\rho = 10^{-5}$  many of the oscillations are gone, and when  $\rho = 10^{-3}$  the computed solution provides a good fit to the exact solution, and the switching points for the discrete problem are roughly within the mesh spacing ( $N = 100$ ) of the exact switching points. In some problems with highly oscillatory solutions, convergence of TV regularized optimal values is established in [13].

When solving (5.1), we also obtain an estimate for the initial costate  $\bar{\mathbf{p}}(0)$  associated with a solution of (1.1). In particular, the KKT conditions at a solution of

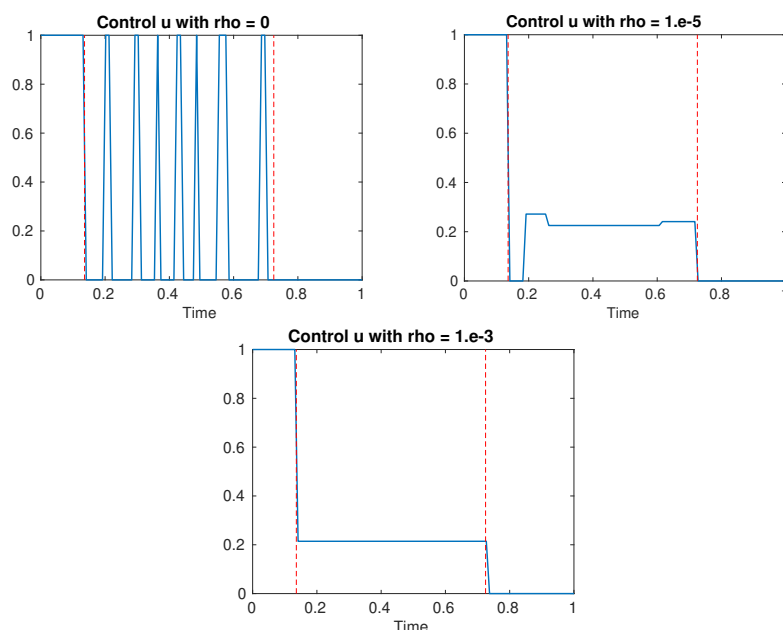


FIG. 5.1. Numerical solutions to (5.1) for the catalyst mixing problem and three different choices for the regularization:  $\rho = 0, 10^{-5}, 10^{-3}$ . Exact switching points appear as dashed lines.

(5.1) yield the following equation for the multiplier  $\mathbf{p}_j$  associated with the constraint  $\mathbf{x}_j + h\mathbf{f}(\mathbf{x}_j, \mathbf{u}_j) - \mathbf{x}_{j+1} = \mathbf{0}$ :

$$\mathbf{p}_{j-1} = \mathbf{p}_j(\mathbf{I} + h\nabla_{\mathbf{x}}\mathbf{f}(\mathbf{x}_j, \mathbf{u}_j)), \quad 1 \leq j \leq N-1, \quad \mathbf{p}_{N-1} = \nabla C(\mathbf{x}_N),$$

where  $\mathbf{p}_0$  is an approximation to  $\bar{\mathbf{p}}(0)$  and  $\mathbf{I}$  is the  $n$  by  $n$  identity matrix.

**6. Numerical studies.** We consider four classic singular control problems from the literature to examine the performance of the switch point algorithm relative to that of previously reported algorithms. The test problems are the following:

1. the catalyst mixing problem, first proposed by Gunn and Thomas [24], and later solved by Jackson [29], with additional analytic formulas given by Li and Liu [31];
2. a problem posed by Jacobson, Gershwin, and Lele [30];
3. a problem posed by Bressan [11];
4. a problem posed by Goddard [23].

All the test problems have known solutions. The switch point algorithm was implemented in MATLAB; the gradients of the objective were evaluated using the formulas given in the paper. The differential equations were integrated using MATLAB ODE45 code, which implements the Dormand–Prince [19] embedded explicit Runge–Kutta (4,5) scheme (both fourth and fifth order accuracy with error estimation). The optimization was performed using PASA [27] (available from Hager’s web page). The experiments were performed on a Dell T7610 workstation with 3.40 GHz processors. We did not implement other algorithms; we simply compare to previously reported results in the literature.

**6.1. Catalyst mixing problem.** The optimal control problem is as follows:

$$\begin{aligned}
 (6.1) \quad & \min a(T) + b(T) - 1 \\
 & \text{s.t. } \dot{a}(t) = -u(t)(k_1 a(t) - k_2 b(t)), \\
 & \dot{b}(t) = u(t)(k_1 a(t) - k_2 b(t)) - (1 - u(t))k_3 b(t), \\
 & a(0) = 1, \quad b(0) = 0, \quad 0 \leq u(t) \leq 1.
 \end{aligned}$$

Here  $a$  and  $b$  are the mole fractions of substances  $A$  and  $B$  which are catalyzed by fraction  $u$  to produce  $C$  in a reactor of length  $T$ . The parameters  $k_i$ ,  $i = 1, 2, 3$ , are given constants, and the objective corresponds to the maximization of  $C$ . Symbolically, the relation is denoted  $A \rightleftharpoons B \rightarrow C$ . As seen in Figure 5.1, the pattern of the optimal control is bang, singular, and off.

The Hamiltonian for this problem is

$$H(a, b, u, p_1, p_2) = [(p_2 - p_1)(k_1 a - k_2 b) + p_2 k_3 b]u - p_2 k_3 b.$$

The switching function, which corresponds to the coefficient of  $u$ , is given by

$$(6.2) \quad \mathcal{S}(t) = (p_2(t) - p_1(t))(k_1 a(t) - k_2 b(t)) + k_3 p_2(t)b(t),$$

where  $\mathbf{p}$  is the solution of the system

$$\begin{aligned}
 \dot{p}_1(t) &= -(p_2(t) - p_1(t))k_1 u(t), & p_1(T) &= 1, \\
 \dot{p}_2(t) &= (p_2(t) - p_1(t))k_2 u(t) + k_3(1 - u(t))p_2(t), & p_2(T) &= 1.
 \end{aligned}$$

When  $\mathcal{S}(t) < 0$ ,  $u(t) = 1$ ; when  $\mathcal{S}(t) > 0$ ,  $u(t) = 0$ ; and when  $\mathcal{S}(t) = 0$ ,  $u$  is singular. In the singular region, both  $\mathcal{S}$  and its derivatives vanish. Differentiating the switch function, we have

$$\begin{aligned}
 \dot{\mathcal{S}}(t) &= k_3 [k_1 a(t)p_2(t) - k_2 b(t)p_1(t)], \\
 \ddot{\mathcal{S}}(t) &= k_3 \left\{ u(t)p_1(t)[k_2 b(t)(k_2 - k_3 - k_1) - 2k_1 k_2 a(t)] \right. \\
 &\quad \left. + u(t)p_2(t)[k_1 a(t)(k_2 - k_3 - k_1) + 2k_1 k_2 b(t)] \right. \\
 &\quad \left. + k_3 [k_1 p_2(t)a(t) + k_2 p_1(t)b(t)] \right\}.
 \end{aligned}$$

Although the control is missing from  $\dot{\mathcal{S}}$ , it appears in  $\ddot{\mathcal{S}}$ . Hence, we use the equation  $\ddot{\mathcal{S}}(t) = 0$  to solve for the control in the singular region. In the following equation, the “(t)” arguments for the state and costate are omitted:

$$(6.3) \quad u_{\text{sing}} = \frac{-k_3(k_1 a p_2 + k_2 b p_1)}{p_1[k_2 b(k_2 - k_3 - k_1) - 2k_1 k_2 a] + p_2[k_1 a(k_2 - k_3 - k_1) + 2k_1 k_2 b]}.$$

With further analysis, it can be shown that the singular control is constant. Jackson [29] derives the following expression for the singular control, where the numeric value given below corresponds to the parameter values  $k_1 = k_3 = 1$  and  $k_2 = 10$  which are used throughout the literature:

$$(6.4) \quad u_{\text{sing}} = \frac{\alpha(1 + \alpha)}{\beta + (1 + \alpha)^2} \approx 0.227142082708498,$$

where  $\alpha = \sqrt{k_3/k_2}$  and  $\beta = k_1/k_2$ .



There are two switch points for the catalyst mixing problem. Analytic formulas for the switching times, presented by Jackson, are

$$s_1 = \left( \frac{1}{k_2(1+\beta)} \right) \log \left( \frac{1+\alpha+\beta}{\alpha} \right) \approx 0.136299034594555,$$

$$s_2 = T - k_3^{-1} \log(1+\alpha) \approx T - 0.274769892408345.$$

The optimal control is

$$u(t) = \begin{cases} 1 & \text{if } 0 \leq t < s_1, \\ u_{\text{sing}} & \text{if } s_1 \leq t < s_2, \\ 0 & \text{if } s_2 \leq t \leq T. \end{cases}$$

An analytic formula is given for the optimal objective value in [31]; the numeric values for the optimal objectives are

$$\begin{aligned} -0.048055685860877 & \quad (T = 1), \\ -0.191814356325161 & \quad (T = 4), \\ -0.477712020050041 & \quad (T = 12). \end{aligned}$$

We solve the test problem using both the Case 1 representation in (6.4), where the exact form of the singular control is exploited, and the Case 2 representation for the singular control given in (6.3) where the algorithm computes the control in the singular region. The versions of the switch point algorithm corresponding to these two representation of the singular control are denoted SPA1 (Case 1) and SPA2 (Case 2).

In Table 6.1 we compare the performance of SPA1 and SPA2 to results from the literature for the problem (6.1) with reactor lengths  $T = 1, 4$ , and  $12$ . The accuracy tolerances used for PASA and ODE45 were  $10^{-8}$ . The starting guesses for the

TABLE 6.1  
Performance and absolute errors for the catalyst mixing problem.

Method	$T$	CPU (s)	$C$ error	$s_1$ error	$s_2$ error
SPA1	1	0.10	$1.6 \times 10^{-10}$	$3.1 \times 10^{-09}$	$1.2 \times 10^{-11}$
SPA2	1	0.30	$9.6 \times 10^{-12}$	$1.9 \times 10^{-10}$	$6.2 \times 10^{-11}$
[7]	1	40 – 100	$5.7 \times 10^{-06}$	—	—
[8]	1	9.74	$2.4 \times 10^{-05}$	—	—
[16]	1	—	$5.7 \times 10^{-06}$	—	—
[22]	1	—	$6.3 \times 10^{-09}$	$3.0 \times 10^{-09}$	$1.1 \times 10^{-12}$
[28]	1	—	$5.7 \times 10^{-06}$	—	—
[31]	1	—	$5.9 \times 10^{-09}$	$1.5 \times 10^{-08}$	$6.8 \times 10^{-08}$
[34]	1	17.90	$1.4 \times 10^{-08}$	—	—
IDE [35]	1	—	$2.3 \times 10^{-05}$	$8.3 \times 10^{-03}$	$7.8 \times 10^{-03}$
[47]	1	90.00	$1.4 \times 10^{-05}$	—	—
[48]	1	38.10	$1.4 \times 10^{-08}$	$2.1 \times 10^{-05}$	$4.9 \times 10^{-04}$
SPA1	4	0.12	$1.1 \times 10^{-10}$	$4.5 \times 10^{-09}$	$1.5 \times 10^{-09}$
SPA2	4	0.68	$1.4 \times 10^{-10}$	$2.4 \times 10^{-10}$	$1.5 \times 10^{-11}$
[4]	4	0.90	—	$4.6 \times 10^{-09}$	$7.6 \times 10^{-09}$
[5]	4	0.33	—	$3.3 \times 10^{-03}$	$4.8 \times 10^{-03}$
[14]	4	1.35	—	$8.4 \times 10^{-08}$	$3.2 \times 10^{-07}$
[15]	4	0.90	—	$5.0 \times 10^{-07}$	$8.3 \times 10^{-06}$
SPA1	12	0.19	$1.7 \times 10^{-10}$	$3.7 \times 10^{-10}$	$4.4 \times 10^{-08}$
SPA2	12	0.98	$2.0 \times 10^{-11}$	$1.6 \times 10^{-09}$	$3.6 \times 10^{-14}$
[16]	12	595.00	$7.7 \times 10^{-04}$	—	—
[31]	12	—	$5.0 \times 10^{-11}$	$1.1 \times 10^{-07}$	$4.0 \times 10^{-07}$
[34]	12	17.79	$7.7 \times 10^{-04}$	—	—
[41]	12	—	$2.7 \times 10^{-04}$	—	—
[43]	12	0.24	$1.6 \times 10^{-03}$	—	—

switching times and the initial costate were accurate to roughly one significant digit. For example, with  $T = 1$  the starting guesses were  $s_1 = 0.1$ ,  $s_2 = 0.7$ ,  $p_1(0) = 0.9$ , and  $p_2(0) = 0.8$ . As seen in Table 6.1, the accuracy was improved from the initial one significant digit to between 9 and 11 significant digits when using  $10^{-8}$  tolerances for both PASA and ODE45.

Note that many of the algorithms in Table 6.1 exploit the known form of the singular control; the computing times for SPA1 where the known form of the singular control is exploited are significantly smaller than the time for SPA2 where the singular control is computed. It is difficult to compare the computing times in Table 6.1 since the computers used to solve the test problem vary widely in speed. Moreover, it should be possible to significantly lower the computing time for the switch point algorithm by developing an implementation in a compiled language instead of MATLAB. And with an ODE integrator tailored to the structure of the control problem, the computing time could be reduced further. Observe that the accuracy of the solution computed by switch point algorithm was relatively high when using a modest  $10^{-8}$  accuracy tolerance for both the optimizer and the ODE integrator.

**6.2. Jacobson's problem [30].** The test problem is given by

$$\begin{aligned} \min \quad & \frac{1}{2} \int_0^5 x_1^2(t) + x_2^2(t) dt \\ \text{s.t.} \quad & \dot{x}_1(t) = x_2(t), \quad \dot{x}_2(t) = u(t), \\ & x_1(0) = 0, \quad x_2(0) = 1, \quad -1 \leq u(t) \leq 1. \end{aligned}$$

This problem as well as the next can be reformulated in form (1.1) by adding a new state variable whose dynamics is the integrand of the objective. After this reformulation, one finds that the costate associated with the new variable is  $p_0 := 1$ , so the Hamiltonian simplifies to

$$H(\mathbf{x}, u, \mathbf{p}) = \frac{1}{2}(x_1^2 + x_2^2) + p_1 x_2 + p_2 u.$$

The switching function is  $\mathcal{S}(t) = p_2(t)$ , where  $\mathbf{p}$  is the solution of the system

$$\begin{aligned} \dot{p}_1(t) &= -x_1(t), & p_1(5) &= 0, \\ \dot{p}_2(t) &= -p_1(t) - x_2(t), & p_2(5) &= 0. \end{aligned}$$

The first two derivatives of the switching function are  $\dot{\mathcal{S}}(t) = -p_1(t) - x_2(t)$  and  $\ddot{\mathcal{S}}(t) = x_1(t) - u(t)$ . In the singular region,  $\ddot{\mathcal{S}} = 0$ , which implies that  $u(t) = x_1(t)$ . The optimal control has one switching point whose first few digits are

$$s_1 \approx 1.41376408763006415924,$$

which is the root of the equation

$$1 - s^2/2 = e^{2s-10}(-1 + 2s - s^2/2).$$

The optimal control is

$$u(t) = \begin{cases} -1 & \text{if } 0 \leq t < s_1, \\ x_1(t) & \text{if } s_1 < t \leq 5, \end{cases}$$

where  $[s_1, 5]$  is the singular interval.

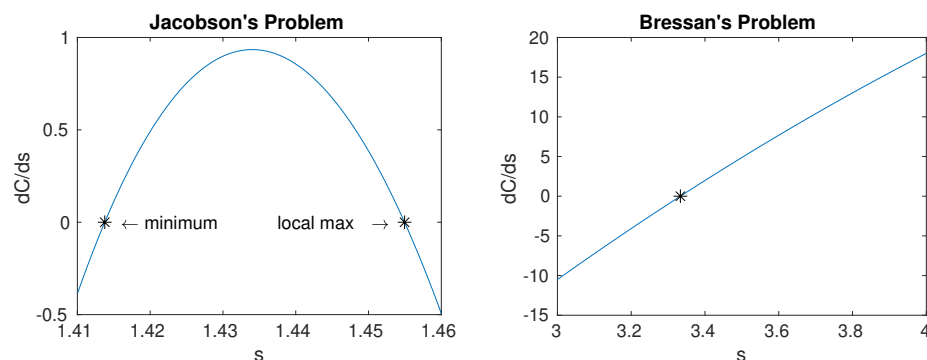


FIG. 6.1. Derivative of the objective versus switch point (zeros marked by stars).

TABLE 6.2  
Performance and absolute errors for Jacobson's problem.

Method	CPU (s)	$s_1$ error
SPA1	0.15	$5.0 \times 10^{-11}$
[4]	2.34	$6.8 \times 10^{-08}$
[5]	0.39	$3.8 \times 10^{-03}$
[14]	132.03	$7.4 \times 10^{-06}$
[15]	0.33	$3.2 \times 10^{-07}$

Since this problem has only one switch point, we will compute it by using the secant method to find a zero of the objective derivative with respect to  $s_1$ . The left panel of Figure 6.1 plots the derivative of the objective function with respect to the switch point. Notice that the derivative vanishes twice, once on the interval  $[1.41, 1.42]$  (corresponding to the optimal control) and once on the interval  $[1.45, 1.46]$  (corresponding to a local maximum). Thus to compute the correct switching point using the secant method, we should start to the left of the local maximum in Figure 6.1. The results in Table 6.2 correspond to the starting iterates 1.42 and 1.41, which bracket the switch point of the optimal control. In five iterations of the secant method, we obtain roughly 11 digit accuracy. The solution time of the switch point algorithm is basically the time of the MATLAB ODE45 integrator. In Figure 6.1 it appears that by choosing the switch point to the right of the local maximum, it may be possible to achieve a smaller value for the cost function. However, for the switch point corresponding to the local maximum in Figure 6.1, the singular control  $u(t) = x_1(t)$  is  $-1$  at  $t = 5$ , and as the switching point moves further to the right, the singular control becomes infeasible with  $u(5) < -1$ .

**6.3. Bressan's problem [11].** The test problem is

$$\begin{aligned}
 \min \quad & \int_0^T x_1^2(t) - x_2(t) \, dt \\
 \text{s.t.} \quad & \dot{x}_1(t) = u(t), \quad \dot{x}_2(t) = -x_1(t), \\
 & x_1(0) = 0, \quad x_2(0) = 0, \quad x_3(0) = 0, \quad -1 \leq u(t) \leq 1,
 \end{aligned}$$

where  $T = 10$  in [5, 14]. The Hamiltonian is

$$H(\mathbf{x}, u, \mathbf{p}) = x_1^2 - x_2 + p_1 u - p_2 x_1.$$

The switching function is  $\mathcal{S}(t) = p_1(t)$  where  $\mathbf{p}$  is the solution of the system

$$\begin{aligned} \dot{p}_1(t) &= p_2(t) - 2x_1(t), & p_1(T) &= 0, \\ \dot{p}_2(t) &= 1, & p_2(T) &= 0. \end{aligned}$$

The first two derivatives of the switching function are  $\dot{\mathcal{S}}(t) = p_2(t) - 2x_1(t)$  and  $\ddot{\mathcal{S}}(t) = 1 - 2u(t)$ . In the singular region,  $\ddot{\mathcal{S}} = 0$ , which implies that  $u(t) = 1/2$ . It can be shown [11] that there is one switching point  $s_1 = T/3$ . The optimal control is

$$u(t) = \begin{cases} -1 & \text{if } 0 \leq t < s_1, \\ 1/2 & \text{if } s_1 < t \leq T. \end{cases}$$

Unlike Jacobson's problem, the plot of the objective derivative versus the switch point (right panel of Figure 6.1) is nearly linear in a wide interval around the switch point  $10/3$  when  $T = 10$ . Starting from the initial iterates 3.0 and 4.0, the secant method converges to 15 digit accuracy in five iterations.

**6.4. Goddard's problem [23].** A classic problem with a free terminal time is Goddard's rocket problem. It is difficult to compare to other algorithms in the literature since there are many variations of Goddard's problem [9, 10, 12, 21, 20, 36, 37, 44, 50], and different algorithms are tested using different versions of the problem. The formulation of the Goddard problem that we use is based on parameters and constraints from both [9, p. 213] and [44, Ex. 3]:

$$\begin{aligned} \min \quad & -h(T) \\ \text{s.t.} \quad & \dot{h}(t) = v(t), & h(0) &= 0, \\ & \dot{v}(t) = \frac{1}{m} [u(t) - \sigma v^2 e^{-h(t)/h_0} - g], & v(0) &= 0, \\ & \dot{m}(t) = -u(t)/c, & m(0) &= 3, \\ & m(t) \geq 1, \quad 0 \leq u(t) \leq u_{\max}, \quad t \in [0, T], \end{aligned}$$

where  $u_{\max} = 193$ ,  $g = 32.174$ ,  $\sigma = 5.4915 \times 10^{-5}$ ,  $c = 1580.9425$ , and  $h_0 = 23800$ . In this problem,  $h$  is the height of a rocket,  $v$  is its velocity,  $m$  is the mass,  $c$  is the exhaust velocity of the propellant,  $g$  is the gravitational acceleration,  $h_0$  is the density scale height,  $u$  is the thrust, and the terminal time is free. The goal is to choose  $T$  and the thrust  $u$  so as to achieve the greatest possible terminal height for the rocket. The mass has a lower bound (the weight of the rocket minus the weight of the propellant) which is taken to be one.

To solve the Goddard rocket problem with the switch point algorithm, we will convert the state constraint  $m \geq 1$  into an additional cost term in the objective. Observe that the mass is a monotone decreasing function of time since  $u \geq 0$ . Since the goal is to achieve the greatest possible height, all the fuel will be consumed and  $m(T) = 1$ . Consequently, the state constraint is satisfied when the terminal constraint  $m(T) = 1$  is fulfilled. By adding a term of the form  $\beta(m(T) - 1)$  to the objective, where  $\beta$  corresponds to the optimal costate evaluated at the terminal time, the solution of the Goddard problem becomes a stationary point of the problem with the modified objective and with the state constraint omitted. To achieve an objective where the solution to the Goddard problem is a local minimizer, we need to incorporate a penalty term in the objective:

$$(6.5) \quad -h(T) + \beta(m(T) - 1) + \frac{\rho}{2}(m(T) - 1)^2$$

TABLE 6.3  
Performance and absolute errors for Bressan's problem.

Method	CPU (s)	$s_1$ error
SPA1	0.09	$1.8 \times 10^{-15}$
[5]	0.29	$3.3 \times 10^{-03}$
[14]	24.80	$5.9 \times 10^{-07}$

with  $\rho > 0$  and  $\beta \approx -2.31774080357308 \times 10^4$ . Our Goddard test problem corresponds to the original Goddard problem, but with the state constraint dropped and with the objective replaced by (6.5).

The optimal solution of the Goddard rocket problem has two switch points and the optimal control is

$$(6.6) \quad u(t) = \begin{cases} u_{\max}, & 0 \leq t \leq s_1 \approx 13.75532627577406, \\ u_{\text{sing}}(t), & s_1 < t \leq s_2 \approx 21.98890645593362, \\ 0, & s_2 < t \leq T \approx 42.88910958027504, \end{cases}$$

where the control in the singular region, gotten from the second derivative of the switching function, is

$$u_{\text{sing}}(t) = \sigma v^2(t) e^{h(t)/h_0} + mg + \frac{mg}{1 + 4\kappa(t) + 2\kappa^2(t)} \left[ \frac{c^2}{h_0 g} (1 + \kappa^{-1}(t)) - 1 - 2\kappa(t) \right],$$

where  $\kappa(t) = c/v(t)$ . The switching points in (6.6) were estimated by integrating forward the dynamics utilizing the known structure for the optimal solution.

We solve the Goddard test problem using  $\rho = 10^5$ , the starting guess  $s_1 = 13$ ,  $s_2 = 21$ , and  $T = 42$ , and the optimizer PASA with the MATLAB integrator ODE45 and computing tolerances  $10^{-8}$ . The solution time was 1.21 s, and the absolute errors in  $s_1$ ,  $s_2$ , and  $T$  were  $1.3 \times 10^{-8}$ ,  $6.0 \times 10^{-8}$ , and  $9.4 \times 10^{-8}$ , respectively. PASA achieved the convergence tolerance in 11 iterations, and essentially all the computing time is associated with ODE45.

**7. Conclusions.** A new approach, the switch point algorithm, was introduced for solving nonsmooth optimal control problems whose solutions are bang-bang, singular, or both bang-bang and singular. For problems where the optimal control has the form  $\mathbf{u}(t) = \phi_j(\mathbf{x}(t), t)$  (Case 1) or  $\mathbf{u}(t) = \phi_j(\mathbf{x}(t), \mathbf{p}(t), t)$  (Case 2) for  $t \in (s_j, s_{j+1})$ ,  $0 \leq j \leq k$ , with  $\mathbf{u}(t)$  feasible for  $\mathbf{s}$  in a neighborhood of the optimal switching points and for the initial costate  $\mathbf{p}(0)$  in a neighborhood of the associated optimal costate, the solution of the optimal control problem reduces to an optimization over the switching points and the choice of the initial costate  $\mathbf{p}(0)$ . Formulas were obtained for the derivatives of the objective with respect to the switch points, the initial costate  $\mathbf{p}(0)$ , and the terminal time  $T$ . All these derivatives can be computed simultaneously with just one integration of the state or generalized state dynamics, followed by one integration of the costate or generalized costate dynamics. A series of test problems were solved by either optimizing over the switching points (and over  $\mathbf{p}(0)$  in Case 2) or computing a point where the derivative of the objective with respect to a switch point vanishes. Accurate solutions were obtained relatively quickly as seen in the comparisons given in Tables 6.1–6.3.

## REFERENCES

- [1] A. A. AGRACHEV, G. STEFANI, AND P. ZEZZA, *Strong optimality for a bang-bang trajectory*, SIAM J. Control Optim., 41 (2002), pp. 991–1014.
- [2] G. ALY, *The computation of optimal singular control*, Internat. J. Control, 28 (1978), pp. 681–688.
- [3] G. M. ANDERSON, *An indirect numerical method for the solution of a class of optimal control problems with singular arcs*, IEEE Trans. Automat. Control, (1972), pp. 363–365.
- [4] O. ANDRÉS-MARTÍNEZ, L. T. BIEGLER, AND A. FLORES-TLACUAHUAC, *An indirect approach for singular optimal control problems*, Computers Chem. Eng., 139 (2020), 106923.
- [5] O. ANDRÉS-MARTÍNEZ, A. FLORES-TLACUAHUAC, S. KAMESWARAN, AND L. T. BIEGLER, *An efficient direct/indirect transcription approach for singular optimal control*, Amer. Inst. Chemical Engineers J., 65 (2019), pp. 937–946.
- [6] S. ATKINS, M. AGHAEE, M. MARTCHEVA, AND W. HAGER, *Solving Singular Control Problems in Mathematical Biology, Using PASA*, arXiv:2010.06744, 2020.
- [7] J. BANGA, R. IRIZARRY-RIVERA, AND W. D. SEIDER, *Stochastic optimization for optimal and model-predictive control*, Computers Chem. Eng., 22 (1998), pp. 603–612.
- [8] M. BELL AND R. SARGENT, *Optimal control of inequality constrained DAE systems*, Computers Chem. Eng., 24 (2000), pp. 2385–2404.
- [9] J. T. BETTS, *Practical Methods for Optimal Control Using Nonlinear Programming*, 2nd ed., SIAM, Philadelphia, 2010.
- [10] J. F. BONNANS, P. MARTINON, AND E. TRÉLAT, *Singular arcs in the generalized Goddard's problem*, J. Optim. Theory Appl., 139 (2008), pp. 439–461.
- [11] A. BRESSAN, *Viscosity Solutions of Hamilton-Jacobi Equations and Optimal Control Problems*, Tech. report, Penn State University, State College, PA, 2011, <https://www.math.psu.edu/bressan/PSPDF/HJ-lnotes.pdf>.
- [12] A. E. BRYSON AND Y.-C. HO, *Applied Optimal Control*, Hemisphere Publishing, New York, 1975.
- [13] M. CAPONIGRO, R. GHEZZI, B. PICCOLI, AND E. TRÉLAT, *Regularization of chattering phenomena via bounded variation controls*, IEEE Trans. Automat. Control, 63 (2018), pp. 2046–2060.
- [14] W. CHEN AND L. T. BIEGLER, *Nested direct transcription optimization for singular optimal control problems*, Amer. Inst. Chemical Engineers J., 62 (2016), pp. 3611–3627.
- [15] W. CHEN, Y. REN, G. ZHANG, AND L. T. BIEGLER, *A simultaneous approach for singular optimal control based on partial moving grid*, Amer. Inst. Chemical Engineers J., 65 (2019), e16584.
- [16] S. DADEBO AND K. MCAULEY, *Dynamic optimization of constrained chemical engineering problems using dynamic programming*, Computers Chem. Eng., 19 (1995), pp. 513–525.
- [17] C. L. DARBY, W. W. HAGER, AND A. V. RAO, *Direct trajectory optimization using a variable low-order adaptive pseudospectral method*, J. Spacecraft Rockets, 48 (2011), pp. 433–445.
- [18] C. L. DARBY, W. W. HAGER, AND A. V. RAO, *An hp-adaptive pseudospectral method for solving optimal control problems*, Optim. Control Appl. Meth., 32 (2011), pp. 476–502, <https://doi.org/10.1002/oca.957>.
- [19] J. R. DORMAND AND P. J. PRINCE, *A family of embedded Runge-Kutta formulae*, J. Comput. Appl. Math., 6 (1980), pp. 19–23.
- [20] Z. FOROOZANDEH, M. SHAMSI, V. AZHMYAKOV, AND M. SHAFIEE, *A modified pseudospectral method for solving trajectory optimization problems with singular arc*, Math. Methods Appl. Sci., 40 (2017), pp. 1783–1793.
- [21] Z. FOROOZANDEH, M. SHAMSI, AND M. D. R. DE PINHO, *A mixed-binary non-linear programming approach for the numerical solution of a family of singular control problems*, Internat. J. Control, 92 (2019), pp. 1551–1566, <https://doi.org/10.1080/00207179.2017.1399216>.
- [22] Z. FOROOZANDEH, M. SHAMSI, AND M. D. R. DE PINHO, *A hybrid direct-indirect approach for solving the singular optimal control problems of finite and infinite order*, Iran. J. Sci. Technol. Trans. Sci., 42 (2018), pp. 1545–1554.
- [23] R. H. GODDARD, *A method of reaching extreme altitudes*, Smithsonian Institution Misc. Collection, 71 (1919).
- [24] D. J. GUNN AND W. J. THOMAS, *Mass transport and chemical reaction in multifunctional catalyst systems*, Chem. Eng. Sci., 20 (1965), pp. 89–100.
- [25] W. W. HAGER, H. HOU, S. MOHAPATRA, A. V. RAO, AND X.-S. WANG, *Convergence rate for a Radau hp-collocation method applied to constrained optimal control*, Comput. Optim. Appl., 74 (2019), pp. 274–314.

- [26] W. W. HAGER AND R. ROSTAMIAN, *Optimal coatings, bang-bang controls, and gradient techniques*, Optim. Control Appl. Methods, 8 (1987), pp. 1–20.
- [27] W. W. HAGER AND H. ZHANG, *An active set algorithm for nonlinear optimization with polyhedral constraints*, Sci. China Math., 59 (2016), pp. 1525–1542.
- [28] R. IRIZARRY, *A generalized framework for solving dynamic optimization problems using the artificial chemical process paradigm: Applications to particulate processes and discrete dynamic systems*, Chem. Eng. Sci., 60 (2005), pp. 5663–5681.
- [29] R. JACKSON, *Optimal use of mixed catalysts for two successive chemical reactions*, J. Optim. Theory Appl., 2 (1968), pp. 27–39.
- [30] D. H. JACOBON, S. B. GERSHWIN, AND M. M. LELE, *Computation of optimal singular controls*, IEEE Trans. Automat. Control, 15 (1970), pp. 67–73.
- [31] G. LI AND X. LIU, *Comments on “optimal use of mixed catalysts for two successive chemical reactions,”* J. Optim. Theory Appl., 165 (2015), pp. 678–692.
- [32] F. LIU, W. W. HAGER, AND A. V. RAO, *Adaptive mesh refinement method for optimal control using nonsmoothness detection and mesh size reduction*, J. Franklin Inst., 352 (2015), pp. 4081–4106.
- [33] F. LIU, W. W. HAGER, AND A. V. RAO, *Adaptive mesh refinement method for optimal control using decay rates of Legendre polynomial coefficients*, IEEE Trans. Control Sys. Tech., 26 (2018), pp. 1475–1483.
- [34] X. LIU, L. CHEN, AND Y. HU, *Solution of chemical dynamic optimization using the simultaneous strategies*, Chinese J. Chem. Eng., 21 (2013), pp. 55–63.
- [35] F. S. LOBATO, V. STEFFEN, AND A. J. S. NETO, *Solution of singular optimal control problems using the improved differential evolution algorithm*, J. Artificial Intelligence Soft Computing Res., 1 (2011), pp. 195–206.
- [36] P. MARTINON, F. BONNANS, J. LAURENT-VARIN, AND E. TRÉLAT, *Numerical study of optimal trajectories with singular arcs for an Ariane 5 launcher*, J. Guid. Control Dyn., 32 (2009), pp. 51–55.
- [37] H. MAURER, *Numerical solution of singular control problems using multiple shooting techniques*, J. Optim. Theory Appl., 18 (1976), pp. 235–257.
- [38] H. MAURER, C. BÜSKENS, J.-H. R. KIM, AND C. Y. KAYA, *Optimization methods for the verification of second order sufficient conditions for bang–bang controls*, Optim. Control Appl. Methods, 26 (2005), pp. 129–156.
- [39] N. P. OSMOLOVSKII AND H. MAURER, *Equivalence of second order optimality conditions for bang–bang control problems. Part 2: Proofs, variational derivatives and representations*, Control Cybernet., 36 (2007), pp. 5–45.
- [40] M. A. PATTERSON, W. W. HAGER, AND A. V. RAO, *A ph mesh refinement method for optimal control*, Optim. Control Appl. Methods, 36 (2015), pp. 398–421.
- [41] D. PHAM, Q. T. PHAM, A. GHANBARZADEH, AND M. CASTELLANI, *Dynamic optimisation of chemical engineering processes using the bees algorithm*, IFAC Proc. Vol., 41 (2008), pp. 6100–6105.
- [42] L. S. PONTYAGIN, V. G. BOLTYANSKII, R. V. GAMKRELIDZE, AND E. F. MISHCHENKO, *The Mathematical Theory of Optimal Processes*, John Wiley, New York, 1962.
- [43] J. RAJESH, K. GUPTA, H. S. KUSUMAKAR, V. JAYARAMAN, AND B. KULKARNI, *Dynamic optimization of chemical processes using ant colony framework*, Computers Chemistry, 25 (2001), pp. 583–95.
- [44] A. V. RAO, D. A. BENSON, C. DARBY, M. A. PATTERSON, C. FRANÇOLIN, I. SANDERS, AND G. T. HUNTINGTON, *Algorithm 902: GPOPS, A MATLAB software for solving multiple-phase optimal control problems using the Gauss pseudospectral method*, ACM Trans. Math. Software, 37 (2010), pp. 1–39, <https://doi.org/10.1145/1731022.1731032>.
- [45] L. RUDIN, S. OSHER, AND E. FATEMI, *Nonlinear total variation based noise removal algorithms*, Phys. D, 60 (1992), pp. 259–268.
- [46] H. SCHÄTTLER AND U. LEDZEWICZ, *Geometric Optimal Control*, Springer, New York, 2012.
- [47] P. TANARTKIT AND L. BIEGLER, *A nested, simultaneous approach for dynamic optimization problems-II: the outer problem*, Computers Chem. Eng., 21 (1997), pp. 735–741.
- [48] V. VASSILIADIS, *Computational Solution of Dynamic Optimization Problems with General Differential-Algebraic Constraints*, Ph.D. thesis, Department of Chemical Engineering and Chemical Technology, Imperial College, London, 1993.
- [49] G. VOSSEN, *Numerische Lösungsmethoden, hinreichende Optimalitätsbedingungen und Sensitivitätsanalyse für optimale bang-bang und singuläre Steuerungen*, Ph.D. thesis, Universität Münster, Germany, 2006.
- [50] G. VOSSEN, *Switching time optimization for bang-bang and singular controls*, J. Optim. Theory Appl., 144 (2010), pp. 409–429.

- [51] G. VOSSEN, *Switching Time Optimization for Bang-Bang and Singular Controls: Variational Derivatives and Applications*, Tech. rep., RWTH Aachen University, Germany, 2010, <https://www.hs-niederrhein.de/maschinenbau-verfahrenstechnik/fachbereich/personenseite-vossen>.
- [52] A. WÄCHTER AND L. T. BIEGLER, *On the implementation of a primal-dual interior point filter line search algorithm for large-scale nonlinear programming*, Math. Program., 106 (2006), pp. 25–57.