# MULTIPLIER METHODS FOR NONLINEAR OPTIMAL CONTROL*

WILLIAM W. HAGER†

**Abstract.** Error estimates are derived for an augmented Lagrangian approximation to an optimal control problem. Convex control constraints are treated explicitly while a Lagrange multiplier is introduced for the nonlinear system dynamics. The nonlinear optimal control problem does not fit the classical theory for estimating the error in multiplier approximations, since the natural coercivity assumption is formulated in a Hilbert space where the cost is not differentiable. This discrepancy between the function space setting needed for coercivity and that needed for differentiability is compensated for by regularity results associated with the necessary conditions. The paper concludes with an analysis of the optimal penalty parameter corresponding to a given finite-element discretization.

**Key words.** multiplier method, optimal control, penalty method, finite elements

**AMS(MOS) subject classification.** 49D29

**1. Introduction.** This paper studies augmented Lagrangian approximations to optimal control problems. Given an interval $I = [0, 1]$, the state $x$ is a map from $I$ to $R^n$ while the control $u$ is a map from $I$ to $R^m$. In our model problem, the system dynamics is possibly nonlinear, the cost is possibly nonquadratic, while the control constraint is convex:

$$
(1) \quad
\begin{aligned}
&\text{minimize} \int_I h(x(t), u(t)) \, dt \\
&\text{subject to } \dot{x}(t) = f(x(t), u(t)) \quad \text{and} \quad u(t) \in \Omega \quad \text{a.e. } t \in I, \\
&\quad x(0) = a, \qquad (x, u) \in W^{1,\infty} \times L^\infty,
\end{aligned}
$$

where $f: R^{n+m} \to R^n$, $h: R^{n+m} \to R$, $\Omega \subset R^m$ is nonempty, closed and convex, $a$ is the given starting condition, $L^\infty$ is the space of essentially bounded functions, and $W^{1,\infty}$ is the space of Lipschitz continuous functions (or, equivalently, the space of essentially bounded functions with essentially bounded derivatives).

We consider an augmented Lagrangian approximation to (1) with quadratic penalty:

$$
(2) \quad
\begin{aligned}
&\text{minimize } C(z) + \langle p^h, F(z) \rangle + \frac{1}{2\varepsilon} \langle F(z), F(z) \rangle \\
&\text{subject to } u(t) \in \Omega \quad \text{a.e. } t \in I, \quad x(0) = a, \quad z = (x, u) \in W^{1,\infty} \times L^\infty,
\end{aligned}
$$

where $z$ denotes the pair $(x, u)$, $p^h$ is any approximation to a Lagrange multiplier $p^*$ associated with the differential equation in (1), $\varepsilon > 0$ is the penalty parameter, $C(z)$ is the integral cost in (1), $F(z) = f(z) - \dot{x}$, and $\langle \cdot, \cdot \rangle$ is the $L^2$ inner product associated with square integrable functions. Multiplier methods seem to originate from work by Arrow and Solow [1], Hestenes [19], and Powell [28]. Additional results are developed in the sequence of papers [31]–[33] by Rockafellar. The book [4] by Bertsekas is a comprehensive reference. Results for problems formulated in a Hilbert space appear

in [17], [25], and [26]. Some interesting applications of augmented Lagrangian tech-
niques to boundary-value problems are developed by Fortin, Glowinski, and their
collaborators in [12].

Multiplier methods can be viewed as a generalization of the penalty method—the
penalty method is obtained from the multiplier method by taking $p^h = 0$. Courant [7]
published the first paper studying penalty techniques for the solution of partial differen-
tial equations. In chronological order, other researchers who have studied penalty
techniques for either partial differential equations or optimal control include Russell
[34], Lions [24], Balakrishnan [3], Babuška [2], King [23], Falk [10], Falk and King
[11], Kikuchi [22], Werner [36], Chen and Mills [5], Reddy [29], Kheshgi and Luskin
[21], Chen et al. [6], and Hager [17].

Let us focus in particular on those papers related to optimal control: Russell [34]
considers state constrained problems; a penalty is introduced for the state constraint
and the convergence of solutions for the penalized problems to a solution of the original
problem is established. Lions [24] considers a quadratic cost problem with a convex
control constraint and with the system dynamics described by a parabolic partial
differential equation. After introducing a penalty for the system dynamics, Lions shows
that the solution to the penalized problem converges strongly to the solution of the
original problem. Balakrishnan also uses a penalty term to handle the system dynamics
in [3]. He gives a detailed analysis of unconstrained, quadratic cost problems with
linear system dynamics. For control constrained problems with nonlinear system
dynamics, he shows that the optimal value associated with the penalized cost function
converges to the optimal value associated with the original problem; in addition, a
maximum principle for the penalized problem is developed, and the limit of the
maximum principle as the penalty parameter tends to zero is analyzed. Chen and Mills
[5] consider an unconstrained quadratic cost problem with linear system dynamics
and with an endpoint constraint. A penalty is introduced for the endpoint constraint,
and it is shown that the deviation between the solution to the penalized problem and
the solution to the original problem is bounded by $O(\varepsilon)$. In [6] Chen et al. examine
unconstrained quadratic cost problems with linear system dynamics. A penalty term
is used to handle the system dynamics. Assuming the penalized problem is solved by
the finite-element method, a condition is formulated that leads to the uncoupling of
the penalization error and the discretization error. In [17] Hager analyzes the error in
finite-element approximations to augmented Lagrangians, and applies the results to
optimal control problems with terminal constraints. Both linear and nonlinear problems
are analyzed, as well as problems for which the system dynamics is described by a
partial differential equation.

In this paper, a penalty is introduced for the system dynamics in a nonlinear
control problem with control constraints. We obtain precise estimates for the distance
between a solution $z_\varepsilon$ to (2) and a solution $z^* = (x^*, u^*)$ to (1). When we try to obtain
error estimates using classical techniques developed for the analysis of multiplier
methods (see [4], [17], or [25]), the following problem is encountered: The cost satisfies
a coercivity assumption in $H^1 \times L^2$ while $f$ and $h$ in (1) are differentiable in the $L^\infty$
norm; here $H^1$ denotes the usual Sobolev space consisting of functions in $L^2$ with an
$L^2$ derivative. This discrepancy between the norm needed for coercivity and the norm
needed for differentiability is compensated for by regularity results associated with the
necessary conditions for the optimal control problem.

The error estimates and analysis in this paper are quite different from the error
estimates contained in our earlier analysis (see [13], [15], [16], and [18]) of multiplier
approximations to convex optimal control problems. In our earlier work, the penalty

term of (2) was not present. If $z^h$ denotes a minimizer associated with the ordinary Lagrangian, then the distance between $z^h$ and $z^*$ in the $L^2$ norm was estimated in terms of the distance between $p^h$ and $p^*$ in the $H^1$ norm. In contrast, the analysis that follows includes a penalty term so that nonconvex problems can be treated; however, our estimate for the distance between $z_\varepsilon$ and $z^*$ is measured in the $H^1 \times L^\infty$ norm using the quantity $\varepsilon \|p^h - p^*\|_{L^2}$. Thus for the augmented Lagrangian, the error only depends on the $L^2$ distance between $p^h$ and $p^*$, not the $H^1$ distance.

Briefly, our paper is organized in the following way: In § 2 we establish regularity results for a linear-quadratic optimal control problem. Later, the nonlinear problem is linearized and the linear-quadratic regularity results are used. Section 3 develops necessary conditions for the nonlinear problem, while § 4 uses the necessary conditions, the regularity results, and the implicit function theorem to estimate the distance between an extreme point $z_\varepsilon$ of (2) and a solution $z^*$ of (1). Using this estimate, we obtain a convergence result for one of the standard iterative implementations of augmented Lagrangian techniques: Letting $p_k$ denote the current approximation to the multiplier $p^*$ associated with the differential equation, the new approximation $z_{k+1}$ to $z^*$ is a local minimizer of the augmented Lagrangian (2) corresponding to $p^h = p_k$; the new approximation $p_{k+1}$ to the multiplier is given by $p_{k+1} = p_k + F(z_{k+1})/\varepsilon$. In § 5, we show that near a local minimizer of (1), extreme points of (2) locally minimize the augmented Lagrangian.

When the augmented Lagrangian (2) is optimized numerically, the space $W^{1,\infty} \times L^\infty$ is replaced by a finite-dimensional approximation. Section 6 examines finite-element approximations. In a simple example, it is seen that, as $\varepsilon \to 0$, while the dimension of the finite-element space is fixed, the finite-dimensional approximation can move away from the solution it is approximating. To achieve good approximation properties as $\varepsilon \to 0$, the dimension of the finite-element space must increase as $\varepsilon$ decreases. We show rigorously that, for a linear quadratic problem with piecewise linear approximations to the state and piecewise constant approximations to the control, the optimal relationship between $\varepsilon$ and the mesh spacing $h$ is $\varepsilon = ch$, where $c$ is an arbitrary positive constant (independent of $h$). For higher-order finite-element spaces, the optimal $\varepsilon$ is bounded by $ch^{2r/3}$, where $r$ is the degree of approximation associated with the finite-element space ($r = 1$ for piecewise linear states and piecewise constant controls).

**2. The linear-quadratic problem.** Our analysis of the augmented Lagrangian (2) is based on properties for a linearization of the original optimal control problem (1). In this section, we study the following linear-quadratic problem:

(3)
$$\text{minimize } \frac{1}{2} \int_I z(t)^T P(t) z(t) \, dt$$
$$\text{subject to } M(z) = 0, \quad u(t) \in \Omega \quad \text{a.e. } t \in I, \quad x(0) = a, \quad z = (x, u) \in H^1 \times L^2,$$

where the superscript $T$ denotes transpose, $P \in L^\infty$, $P(t)$ is a symmetric matrix for almost every $t$, and $M : H^1 \times L^2 \to L^2$ is defined by

$$M(z) = Ax + Bu - \dot{x}$$

with $A \in L^2$ and $B \in L^\infty$. Existence of a solution to (3) is related to the following property of a bilinear form.

LEMMA 1. *Let $\pi$ be a symmetric, continuous bilinear form defined on a nonempty, closed convex subset $K$ of a Hilbert space $V$, and let $\langle \cdot, \cdot \rangle_V$ denote the Hilbert space inner product. If there exists $\alpha > 0$ such that*

(4)
$$\pi(w - v, w - v) \geqq \alpha \langle w - v, w - v \rangle_V \quad \text{for all } w, v \in K,$$

*then for any $\phi \in V$, the quadratic program*

(5)
$$\text{minimize } \tfrac{1}{2} \pi(v, v) - \langle \phi, v \rangle_V$$

$$\text{subject to } v \in K$$

*has a unique solution $w$. This solution is the unique $w \in K$ that satisfies the variational inequality $\pi(w, v - w) \geqq \langle \phi, v - w \rangle_V$ for all $v \in K$. If $w_i$ denotes the solution of (5) corresponding to $\phi = \phi_i$ for $i = 1$ and $i = 2$, then we have*

(6)
$$\alpha \| w_1 - w_2 \| \leqq \| \phi_1 - \phi_2 \|,$$

*where $\| \cdot \|$ is the norm induced by the Hilbert space inner product.*

*Proof.* The first part of the lemma is contained in Chapter 1 of [24]. To establish estimate (6), let us consider the variational inequality $\pi(w, v - w) \geqq \langle \phi, v - w \rangle_V$ for all $v \in K$. Choosing $\phi = \phi_1$, $w = w_1$, and $v = w_2$ followed by $\phi = \phi_2$, $w = w_2$, and $v = w_1$ and adding the resulting relations yields

$$\pi(w_2 - w_1, w_2 - w_1) \leqq \langle \phi_2 - \phi_1, w_2 - w_1 \rangle_V,$$

from which we get (6).    □

There are two different ways to apply Lemma 1 to the linear-quadratic optimal control problem (3). In the first approach, we think of the state and control as independent variables in $H^1 \times L^2$. Thus $\phi = 0$, $v = z$, $\pi(z, z) = \langle z, Pz \rangle$, and $K$ consists of those $z \in H^1 \times L^2$ that satisfy the constraints of (3). In the second approach, we regard the state as an affine function of the control and we take $v = u$. From the system dynamics of (3), $x = L(Bu) + \psi$, where $x = L(y)$ denotes the solution to

$$\dot{x} = Ax + y, \qquad x(0) = 0,$$

and $x = \psi$ denotes the solution to $\dot{x} = Ax$, $x(0) = a$. Partitioning $P$ in the form

(7)
$$P = \begin{bmatrix} Q & S \\ S^T & R \end{bmatrix},$$

where $Q$ is $n \times n$ and $R$ is $m \times m$, the cost function in (3) can be expressed as

$$\langle z, Pz \rangle = \langle L(Bu), QL(Bu) \rangle + 2 \langle L(Bu), Su \rangle + \langle u, Ru \rangle + 2 \langle \psi, Su + QL(Bu) \rangle + \langle \psi, Q\psi \rangle.$$

Hence, we can apply Lemma 1 with $v = u$:

$$\pi(u, u) = \langle L(Bu), QL(Bu) \rangle + 2 \langle L(Bu), Su \rangle + \langle u, Ru \rangle,$$

$$\phi = -(S^T + B^T L^T Q)\psi,$$

where $L^T$ denotes the adjoint of $L$. In this formulation, $K$ consists of those $u \in L^2$ that satisfy the control constraints of (3). That is, $K = U$, where $U$ is the set of controls $u \in L^2$ with $u(t) \in \Omega$ for almost every $t \in I$.

When we identify the $v$ of Lemma 1 with the $z$ of (3) and we think of the state and the control as independent variables, the analogue of the coercivity assumption (4) is the following: There exists $\alpha > 0$ such that

$$\langle z, Pz \rangle \geqq \alpha [\langle x, x \rangle_{H^1} + \langle u, u \rangle] \text{ for all } z = (x, u) \in H_0^1 \times L^2$$

(8)
$$\text{with } M(z) = 0 \text{ and } u = u_1 - u_2$$

$$\text{for some } u_1 \text{ and } u_2 \in U.$$

Here $H_0^1$ is the subspace of $H^1$ consisting of functions that vanish at $t = 0$ and the inner product $\langle \cdot, \cdot \rangle_{H^1}$ is defined by

$$\langle x, x \rangle_{H^1} = \langle \dot{x}, \dot{x} \rangle + \langle x, x \rangle.$$

Observe that if $z = (x, u) \in H_0^1 \times L^2$ and $M(z) = 0$, then there exists a constant $c$, independent of $z$, such that

$$\langle \dot{x}, \dot{x} \rangle + \langle x, x \rangle \leqq c \langle u, u \rangle.$$

Hence, (8) is equivalent to the following condition: There exists $\alpha > 0$ such that

(9)
$$\langle z, Pz \rangle \geqq \alpha \langle u, u \rangle \text{ for all } z = (x, u) \in H_0^1 \times L^2 \text{ with } M(z) = 0 \text{ and } u = u_1 - u_2$$
for some $u_1$ and $u_2 \in U$.

In other words, the cost needs only to be coercive in the control since coercivity in the state holds automatically. Note that, when we take $v = u$ in Lemma 1 and we view the state as an affine function of the control, assumption (4) applied to (3) reduces to (9).

By Lemma 1, (9) implies that there exists a unique solution $z^* = (x^*, u^*)$ to (3). Moreover, if $z^*$ is any solution of (3), then for any $z$ that is feasible in (3), the convexity of the constraints implies that

$$\frac{d}{d\tau} \{ \langle P(z^* + \tau(z - z^*)), z^* + \tau(z - z^*) \rangle \}_{\tau=0} \geqq 0,$$

which yields

(10)
$$\langle Pz^*, z - z^* \rangle \geqq 0.$$

The necessary conditions for (3) are obtained by simplifying this variational inequality using the adjoint variable $p$ defined by

$$\dot{p} + A^T p + Qx + Su = 0 \quad \text{and} \quad p(1) = 0.$$

Let $p = p^*$ denote the solution of the adjoint equation corresponding to $x = x^*$ and $u = u^*$. Adding to (10) the equality

$$\langle p^*, M(z - z^*) \rangle = 0$$

and integrating by parts, we obtain

$$\langle B^T p^* + S^T x^* + Ru^*, v - u^* \rangle \geqq 0 \quad \text{for all } v \in U.$$

Therefore, $x = x^*$, $u = u^*$, and $p = p^*$ is a solution of the system of relations

(11)
$$\dot{p} + A^T p + Qx + Su = 0, \qquad p(1) = 0,$$
$$\langle B^T p + S^T x + Ru, v - u \rangle \geqq 0 \quad \text{for all } v \in U,$$
$$-\dot{x} + Ax + Bu = 0, \qquad x(0) = a.$$

Conversely, if (9) holds, then any solution of (11) yields the optimal solution to (3). To demonstrate this, suppose that $x = x^*$, $u = u^*$, and $p = p^*$ is a solution of (11). Expanding about $z^*$, we have

$$\langle z, Pz \rangle = \langle z^*, Pz^* \rangle + 2 \langle Pz^*, z - z^* \rangle + \langle z - z^*, P(z - z^*) \rangle.$$

If $z$ is feasible in (3), then $\langle z - z^*, P(z - z^*) \rangle \geqq 0$ by (9). Since the inequality in (11) is equivalent to $\langle Pz^*, z - z^* \rangle \geqq 0$, it follows that $\langle z, Pz \rangle \geqq \langle z^*, Pz^* \rangle$ whenever $z$ is feasible in (3). In summary, we have Lemma 2.

LEMMA 2. *Any solution to* (3) *satisfies the relations* (11) *for some p. Moreover, if* (9) *holds, then* (3) *and* (11) *have a unique solution.*

Under hypothesis (9), we now study how the solution to (11) behaves after a linear perturbation. Let $L^p$ denote the usual space of functions consisting of those $f$ with $|f|^p$ summable. Given $q_i$ and $s_i$ in $L^1$ and $r_i$ in $L^2$ (for $i = 1$ and 2), we consider the following perturbation of (11):

(12)
$$\dot{p} + A^T p + Qx + Su = q_i, \qquad\qquad p(1) = 0,$$
$$\langle B^T p + S^T x + Ru, v - u \rangle \geqq \langle r_i, v - u \rangle \quad \text{for all } v \in U,$$
$$-\dot{x} + Ax + Bu = s_i, \qquad\qquad x(0) = a.$$

Observe that (12) is the necessary condition corresponding to the following linear-quadratic problem:

(13)
minimize $\quad \frac{1}{2}\langle z, Pz \rangle - \langle q_i, x \rangle - \langle r_i, u \rangle$

subject to $\quad M(z) = s_i, \quad u(t) \in \Omega \quad \text{a.e. } t \in I, \quad x(0) = a, \quad z = (x, u) \in H^1 \times L^2.$

Hence, there is a one-to-one correspondence between a solution to (12) and a minimizer of (13).

We apply (6) in Lemma 1 to (13), viewing the state as a linear function of the control. After solving for the state in terms of the control, the perturbation $s_i$, and the starting condition $a$, we see that the perturbations $q_i$, $r_i$, and $s_i$ appear as linear parameters in the optimization process. In particular, using the linear operator $L$ and the function $\psi$ introduced after Lemma 1, yields the linear terms

$$\langle \psi, QL(Bu) + Su \rangle - \langle L(s_i), QL(Bu) + Su \rangle - \langle q_i, L(Bu) \rangle - \langle r_i, u \rangle.$$

If (9) holds, then Lemma 1 gives us the estimate

$$\alpha \|u_1 - u_2\|_{L^2} \leqq \|r_1 - r_2 + (S^T + B^T L^T Q)L(s_1 - s_2) + B^T L^T (q_1 - q_2)\|_{L^2}.$$

If we use the standard representation for $L$ in terms of an integral, it follows that

$$\|u_2 - u_1\|_{L^2} \leqq c[\|q_2 - q_1\|_{L^1} + \|r_2 - r_1\|_{L^2} + \|s_2 - s_1\|_{L^1}],$$

where $x_i$, $u_i$, and $p_i$ is the solution of (12) corresponding to $q_i$, $r_i$, and $s_i$. (Here and elsewhere, $c$ denotes a generic constant.)

Let $W^{1,p}$ denote the Sobolev space consisting of functions in $L^p$ with a derivative in $L^p$. (Of course, the superscript $p$ in the spaces $L^p$ and $W^{1,p}$ is unrelated to the costate variable $p$.) From the state equation, it follows that

$$\|x_2 - x_1\|_{W^{1,1}} \leqq c[\|q_2 - q_1\|_{L^1} + \|r_2 - r_1\|_{L^2} + \|s_2 - s_1\|_{L^1}].$$

Also, by the adjoint equation we have

$$\|p_2 - p_1\|_{W^{1,1}} \leqq c[\|q_2 - q_1\|_{L^1} + \|r_2 - r_1\|_{L^2} + \|s_2 - s_1\|_{L^1}].$$

Collecting results gives us

(14) $\quad \|x_2 - x_1\|_{W^{1,1}} + \|u_2 - u_1\|_{L^2} + \|p_2 - p_1\|_{W^{1,1}} \leqq c[\|q_2 - q_1\|_{L^1} + \|r_2 - r_1\|_{L^2} + \|s_2 - s_1\|_{L^1}].$

For nonlinear optimal control problems, this inequality is too weak since the control lies in $L^2$ and the cost is not differentiable in $L^2$. We now give inequalities in stronger norms under the following assumption: There exists $\bar{\alpha} > 0$ (independent of $t$) such that

(15)
$$u^T R(t)u \geqq \bar{\alpha} u^T u$$

whenever $u = u_1 - u_2$ with $u_1$ and $u_2 \in \Omega$. Let us consider the pointwise version of the inequality in (12):

$$(B(t)^T p_i(t) + S(t)^T x_i(t) + R(t)u_i(t))^T (v - u_i(t)) \geqq r_i(t)^T (v - u_i(t))$$

for all $v \in \Omega$ and for almost every $t \in I$. Applying Lemma 1 to this pointwise inequality, or to the equivalent finite-dimensional quadratic program, we have

$$(16) \qquad |u_2(t) - u_1(t)| \leqq c[|x_2(t) - x_1(t)| + |p_2(t) - p_1(t)| + |r_2(t) - r_1(t)|].$$

Since the $L^\infty$ norm of $x_1 - x_2$ and $p_1 - p_2$ is bounded in terms of the $W^{1,1}$ norm, it follows from (14) and (16) that

$$(17) \qquad \|u_2 - u_1\|_{L^\infty} \leqq c[\|q_2 - q_1\|_{L^1} + \|r_2 - r_1\|_{L^\infty} + \|s_2 - s_1\|_{L^1}].$$

More generally, if $q_i$ and $s_i$ lie in $L^p$, then the first and last equation in (12) coupled with (17) yield $W^{1,p}$ bounds for the state and adjoint variables as stated in Lemma 3.

LEMMA 3. *Suppose that $A$, $Q$, $q_i$, and $s_i$ are elements of $L^p$ and $B$, $S$, $R$, and $r_i$ are elements of $L^\infty$ for $i = 1$ and $i = 2$ and for some $p$ between $1$ and $\infty$. If (9) and (15) hold, then the solution to (12) has the following Lipschitz property:*

$$\|x_2 - x_1\|_{W^{1,p}} + \|u_2 - u_1\|_{L^\infty} + \|p_2 - p_1\|_{W^{1,p}} \leqq c[\|q_2 - q_1\|_{L^p} + \|r_2 - r_1\|_{L^\infty} + \|s_2 - s_1\|_{L^p}].$$

**3. Nonlinear formulations.** Now suppose that the nonlinear problem (1) has a solution $z^* = (x^*, u^*)$. We assume that there exists a bounded open set $\Sigma \subset R^{n+m}$, where both $f$ and $h$ are twice continuously differentiable, and that there exists $\delta > 0$ such that $z^*(t) \in \Sigma$ and the distance from $z^*(t)$ to the boundary of $\Sigma$ is at least $\delta$ for almost every $t \in I$. Letting $H$ be the Hamiltonian defined by $H(x, u, p) = h(x, u) + p^T f(x, u)$, the adjoint system associated with (1) is given by

$$\dot{p}(t) = -\frac{\partial H(x(t), u(t), p(t))}{\partial x} \quad \text{a.e. } t \in I, \quad p(1) = 0.$$

If $p = p^*$ denotes the solution to the adjoint equation corresponding to $x = x^*$ and $u = u^*$, it follows from the control minimum principle [27] that

$$\int_I \frac{\partial H(x^*(t), u^*(t), p^*(t))}{\partial u} (v(t) - u^*(t)) \, dt \geqq 0 \quad \text{for all } v \in U.$$

In summary, $x = x^*$, $u = u^*$, and $p = p^*$ satisfy the following necessary conditions:

$$\dot{p} + H_x(x, u, p) = 0, \qquad p(1) = 0,$$

$$(18) \qquad \langle H_u(x, u, p), v - u \rangle \geqq 0 \quad \text{for all } v \in U,$$

$$f(x, u) - \dot{x} = 0, \qquad x(0) = a,$$

where the subscripts $x$ and $u$ denote partial derivatives with respect to $x$ and $u$.

We will show that for $\varepsilon$ sufficiently small, (2) has a local minimizer $z_\varepsilon$ that approaches $z^*$ as $\varepsilon \to 0$. We proceed in the following way. Assuming that there exists a local minimizer $z_\varepsilon$ for (2) that is near $z^*$, we equate the Gateaux derivative to zero to obtain equations for $z_\varepsilon$ and for the multiplier approximation $p_\varepsilon = p^h + F(z_\varepsilon)/\varepsilon$. The equations that we obtain differ from (18) by a small perturbation. Using the Lipschitz result contained in Lemma 3 and the implicit function theorem, we get an estimate for the errors in $z_\varepsilon$ and $p_\varepsilon$. Finally, we show that $z_\varepsilon$ is a local minimizer for the augmented Lagrangian (2).

To begin, suppose that (2) has a local minimizer $z_\varepsilon = (x_\varepsilon, u_\varepsilon)$ near $z^*$. Given $v \in U \cap L^\infty$, we define $\delta u = v - u_\varepsilon$. If $\delta x \in W^{1,\infty}$ with $\delta x(0) = 0$ and $\delta z = (\delta x, \delta u)$, it follows that

$$(19) \quad \frac{d}{d\tau} \left\{ C(z_\varepsilon + \tau \delta z) + \langle p^h, F(z_\varepsilon + \tau \delta z) \rangle + \frac{1}{2\varepsilon} \langle F(z_\varepsilon + \tau \delta z), F(z_\varepsilon + \tau \delta z) \rangle \right\}_{\tau=0} \geqq 0.$$

Taking the derivative in (19), we conclude that

(20) $$\langle h_z(z_\varepsilon), \delta z \rangle + \langle p_\varepsilon, f_z(z_\varepsilon)\delta z - \delta\dot{x} \rangle \geqq 0,$$

where $p_\varepsilon = p^h + F(z_\varepsilon)/\varepsilon$. Since (20) holds for all $\delta x \in W^{1,\infty}$ with $\delta x(0) = 0$, the equivalence class associated with $p_\varepsilon$ contains an absolutely continuous function (see [16]). Hence, there is no loss of generality in assuming that $p_\varepsilon$ is absolutely continuous. Integrating (20) by parts gives

$$-p_\varepsilon(1)\delta x(1) + \langle \dot{p}_\varepsilon + H_x(x_\varepsilon, u_\varepsilon, p_\varepsilon), \delta x \rangle + \langle H_u(x_\varepsilon, u_\varepsilon, p_\varepsilon), \delta u \rangle \geqq 0.$$

Taking $\delta u = 0$ and varying $\delta x$, we conclude that

(21) $$p_\varepsilon(1) = 0 \quad \text{and} \quad \dot{p}_\varepsilon + H_x(x_\varepsilon, u_\varepsilon, p_\varepsilon) = 0.$$

Taking $\delta x = 0$ and replacing $\delta u$ by $v - u_\varepsilon$, we have

$$\langle H_u(x_\varepsilon, u_\varepsilon, p_\varepsilon), v - u_\varepsilon \rangle \geqq 0 \quad \text{for all } v \in U.$$

Finally, let us rearrange the definition $p_\varepsilon = p^h + F(z_\varepsilon)/\varepsilon$ to obtain

$$f(x_\varepsilon, u_\varepsilon) - \dot{x}_\varepsilon - \varepsilon(p_\varepsilon - p^h) = 0.$$

Combining these observations, we see that $x = x_\varepsilon$, $u = u_\varepsilon$, and $p = p_\varepsilon$ satisfy the relations

(22) $$\begin{aligned} \dot{p} + H_x(x, u, p) &= 0, & p(1) &= 0, \\ \langle H_u(x, u, p), v - u \rangle &\geqq 0 & &\text{for all } v \in U, \\ f(x, u) - \dot{x} - \varepsilon(p - p^h) &= 0, & x(0) &= a. \end{aligned}$$

Observe that (18) and (22) are identical except for the $\varepsilon$ term. Hence, $(x^*, u^*, p^*)$ and $(x_\varepsilon, u_\varepsilon, p_\varepsilon)$ satisfy nearly the same equations. We will use the implicit function theorem to estimate the distance between a solution $(x^*, u^*, p^*)$ to (18) and a solution to (22).

**4. Error estimates.** A nice treatment of the implicit function theorem for a variational inequality appears in Robinson's paper [30]. Although Robinson's setting does not exactly fit our application, these differences can be handled with appropriate changes in the problem formulation and in the analysis. For completeness, we now give a development of the implicit function theorem for inequalities, providing explicit estimates for the constants that appear in the bounds and relating the implicit function theorem to the classical contraction mapping principle. Robinson's paper considers an equation involving a parameter, giving estimates for the change in the solution relative to a change in the parameter. The analysis that follows is more in the spirit of our paper [17] in which we estimate the distance between a given point and a root of an equation. In [17] there are no constraints, while in [30] the constraint set is assumed to be convex. In the analysis below, the constraint set is arbitrary.

Let $X$ be a Banach space, let $Y$ be a normed subspace of the dual space $X^*$, and let $K$ be a subset of $X$. Given a map $T$ from $X$ to $Y$, let us consider the following variational inequality.

(23)      Find $\zeta \in K$ such that $\langle T(\zeta), \eta - \zeta \rangle_X \geqq 0$ for all $\eta \in K$.

where $\langle \cdot, \cdot \rangle_X$ denotes the usual pairing between a space and its dual. We will formulate conditions under which (23) has a solution $\zeta_1$ in the neighborhood of some given point $\zeta_0$. Our analysis is based on the classical contraction mapping theorem (for example, see [20, p. 110]), which is stated below.

LEMMA 4. *Suppose that $\Phi$ is a map from the Banach space $X$ to itself, $\theta$ is an element of $X$, $\gamma$ is a real number with $0 \leqq \gamma < 1$, and $r$ is any real number satisfying*

$$(24) \qquad r \geqq \frac{\|\theta - \Phi(\theta)\|}{1 - \gamma}.$$

*If $\Phi$ satisfies the Lipschitz condition*

$$\|\Phi(\zeta) - \Phi(\eta)\| \leqq \gamma \|\zeta - \eta\|$$

*whenever $\zeta$ and $\eta$ lie in $X$ and $\|\zeta - \theta\| \leqq \|\eta - \theta\| \leqq r$, then the fixed-point equation $\zeta = \Phi(\zeta)$ has a unique root inside the ball with center $\theta$ and radius $r$.*

To apply this result to the variational inequality (23), the equation is linearized and a mapping $\Phi$ is constructed to which we apply Lemma 4. Let $L$ be any bounded linear operator mapping $X \to Y$ and consider the following linear variational inequality:

$$(25) \qquad \text{Find } \zeta \in K \text{ such that } \langle L(\zeta - \zeta_0) + \phi, \eta - \zeta \rangle_X \geqq 0 \text{ for all } \eta \in K.$$

It is assumed that for $\phi = y_0 \in Y$, (25) has the solution $\zeta = \zeta_0$. Defining the parameter $D_\rho$ by

$$D_\rho = \sup_{\|\zeta - \zeta_0\| \leqq \rho} \|T'[\zeta] - L\|,$$

where the prime denotes a Fréchet derivative, we have Theorem 1.

THEOREM 1. *Let $\lambda$, $\rho$, and $\sigma$ be parameters that satisfy the relations*

$$0 \leqq \lambda D_\rho < 1, \quad \rho \geqq \frac{\lambda \|T(\zeta_0) - y_0\|}{1 - \lambda D_\rho}, \quad \sigma \geqq \rho D_\rho + \|T(\zeta_0) - y_0\|.$$

*Suppose the $T$ is continuously Fréchet differentiable in a ball with center $\zeta_0$ and radius $\rho$, and for each $\phi$ in $Y$ with $\|\phi - y_0\| \leqq \sigma$, the linearized variational inequality (25) has a solution $\zeta = \Psi(\phi)$ such that $\Psi(y_0) = \zeta_0$ and the Lipschitz condition*

$$\|\Psi(\phi_1) - \Psi(\phi_2)\| \leqq \lambda \|\phi_1 - \phi_2\|$$

*holds whenever $\|\phi_1 - y_0\| \leqq \|\phi_2 - y_0\| \leqq \sigma$. Then variational inequality (23) has a solution $\zeta_1$ that satisfies the inequality*

$$(26) \qquad \|\zeta_1 - \zeta_0\| \leqq \frac{\lambda}{1 - \lambda D_\rho} \|T(\zeta_0) - y_0\|.$$

*Proof.* We apply Lemma 4 with $\theta = \zeta_0$, $r = \lambda \|T(\zeta_0) - y_0\| / (1 - \lambda D_\rho)$, $\gamma = \lambda D_\rho$, and $\Phi(\chi) = \Psi(T(\chi) - L(\chi - \zeta_0))$. Thus a fixed point $\zeta_1$ of $\Phi$ is a solution of the variational inequality (23). To begin, the fundamental theorem of calculus yields

$$T(\zeta) - y_0 - L(\zeta - \zeta_0) = T(\zeta_0) - y_0 + \int_0^1 \{T'[s\zeta + (1-s)\zeta_0] - L\}(\zeta - \zeta_0) \, ds.$$

Taking norms, it follows that $\|T(\zeta) - y_0 - L(\zeta - \zeta_0)\| \leqq \sigma$ whenever $\|\zeta - \zeta_0\| \leqq \rho$. Hence, by the Lipschitz property for $\Psi$, we have

$$\|\Phi(\zeta) - \Phi(\eta)\| \leqq \lambda \|T(\zeta) - T(\eta) - L(\zeta - \eta)\|$$

whenever $\|\zeta - \zeta_0\| \leqq \|\eta - \zeta_0\| \leqq \rho$. Again, expressing $T(\eta) - T(\zeta)$ as the integral of the derivative evaluated along the line segment connecting $\eta$ to $\zeta$, we conclude that

$$\|\Phi(\zeta) - \Phi(\eta)\| \leqq \lambda D_\rho \|\zeta - \eta\|$$

whenever $\|\zeta - \zeta_0\| \leq \|\eta - \zeta_0\| \leq \rho$. Hence, the contraction property of Lemma 4 holds with $\gamma = \lambda D_\rho$. Finally, we estimate the difference $\|\zeta_0 - \Phi(\zeta_0)\|$. Since $\Psi(y_0) = \zeta_0$, it follows that

$$\|\zeta_0 - \Phi(\zeta_0)\| = \|\Psi(y_0) - \Psi(T(\zeta_0))\| \leq \lambda \|y_0 - T(\zeta_0)\|.$$

Condition (24) in Lemma 4 holds since

$$\rho \geq r = \frac{\lambda \|T(\zeta_0) - y_0\|}{1 - \lambda D_\rho} \geq \frac{\|\zeta_0 - \Phi(\zeta_0)\|}{1 - \lambda D_\rho} = \frac{\|\theta - \Phi(\theta)\|}{1 - \gamma}.$$

Hence, by Lemma 4, $\Phi$ has a fixed point $\zeta_1$ and the distance from $\zeta_0$ to $\zeta_1$ is at most $r$. This establishes (26) since the right side of (26) is $r$.      □

*Remark.* Referring to the proof of Lemma 4 and of Theorem 1, we see that $T$ and its derivative are evaluated only in a convex set containing both $\zeta_0$ and points in $K$ near $\zeta_0$.

We apply Theorem 1 to the necessary conditions (18) and (22) associated with the optimal control problem in the following way: The Banach space $X$ corresponding to the triple $(x, u, p)$ is $W^{1,p} \times L^\infty \times W^{1,p}$, where $1 \leq p \leq \infty$, $Y$ is $L^p \times L^\infty \times L^p$, and for $(q, r, s) \in Y$, the associated linear functional that operates on elements $(x, u, p)$ of $W^{1,p} \times L^\infty \times W^{1,p}$ is given by

$$\int_I [q(t)x(t) + r(t)u(t) + s(t)p(t)] \, dt.$$

The constraint set $K$ is the collection of $(x, u, p)$ in $W^{1,p} \times L^\infty \times W^{1,p}$ with $x(0) = a$, $u \in U$, and $p(1) = 0$. The operator $T$ is given by

$$T(x, u, p) = T_0(x, u, p) + \begin{bmatrix} 0 \\ 0 \\ \varepsilon(p^h - p) \end{bmatrix},$$

where

$$T_0(x, u, p) = \begin{bmatrix} \dot{p} + H_x(x, u, p) \\ H_u(x, u, p) \\ f(x, u) - \dot{x} \end{bmatrix}.$$

The point $\zeta_0$ is $(x^*, u^*, p^*)$ while $y_0$ is $T_0(\zeta_0)$ and $L$ is $T_0'(\zeta_0)$. Defining the matrices $A$, $B$, and $P$ by

$$A(t) = \frac{\partial f(x^*(t), u^*(t))}{\partial x}, \quad B(t) = \frac{\partial f(x^*(t), u^*(t))}{\partial u}, \quad P(t) = \frac{\partial^2 H(z^*(t), p^*(t))}{\partial z^2},$$

and partitioning $P$ as in (7), the Fréchet derivative of $T_0$ evaluated at $\zeta_0$ can be expressed as

$$T_0'[\zeta_0](x, u, p) = \begin{bmatrix} \dot{p} + A^T p + Qx + Su \\ B^T p + S^T x + Ru \\ Ax + Bu - \dot{x} \end{bmatrix}.$$

By our differentiability assumptions for $h$ and $f$, the Fréchet derivative of $T$ is continuous in a neighborhood of $\zeta_0$. Under the hypotheses of Lemma 3, linearization (25) has a unique solution for all choices of the perturbation, and this solution is a Lipschitz continuous function of the data. This implies that the $\sigma$ of Theorem 1 is $\infty$. If $\lambda$ denotes the Lipschitz constant of Lemma 3, then we can choose $\rho$ and $\varepsilon$ sufficiently small that

$\lambda D_\rho < 1$. Since $\| T(\zeta_0) - y_0 \| = \varepsilon \| p^h - p^* \|_{L^p}$, we can also choose $\varepsilon$ sufficiently small that $\rho \geqq \lambda \| T(\zeta_0) - y_0 \| / (1 - \lambda D_\rho)$. Hence, the assumptions of Theorem 1 are satisfied, and we have Theorem 2.

THEOREM 2. *Suppose that* (1) *has a solution* $z^*$, *that there exists a bounded open set* $\Sigma \subset R^{n+m}$, *where both $f$ and $h$ are twice continuously differentiable, and that there exists $\delta > 0$ with $z^*(t) \in \Sigma$ and the distance from $z^*(t)$ to the boundary of $\Sigma$ at least $\delta$ for almost every $t \in I$. If* (9) *and* (15) *hold and $p^h \in L^p$ for some $p$ between 1 and $\infty$, then for $\varepsilon \| p^h - p^* \|_{L^p}$ sufficiently small,* (22) *has a solution* $(x, u, p) = (x_\varepsilon, u_\varepsilon, p_\varepsilon)$ *that satisfies*

$$\| x_\varepsilon - x^* \|_{W^{1,p}} + \| u_\varepsilon - u^* \|_{L^\infty} + \| p_\varepsilon - p^* \|_{W^{1,p}} \leqq \frac{\lambda \varepsilon \| p^h - p^* \|_{L^p}}{1 - \lambda D_\rho},$$

*where $D_\rho \to 0$ as $\varepsilon \to 0$ and $\lambda$ is the Lipschitz constant $c$ of Lemma 3.*

In practice, augmented Lagrangians are used in an iterative fashion. At iteration $k$, the current value of the penalty parameter is $\varepsilon_k$, and the current approximation to the multiplier $p^*$ associated with the differential equation is $p_k$. The new approximation $z_{k+1} = (x_{k+1}, u_{k+1})$ to a solution of the optimal control problem (1), and the new approximation $p_{k+1}$ to the multiplier, satisfy (22) when $\varepsilon = \varepsilon_k$ and $p^h = p_k$. Hence, starting from an initial approximation $p_0$ to the multiplier, this iteration generates a sequence $(x_k, u_k, p_k)$, $k = 1, 2, \cdots$, that converges to $(x^*, u^*, p^*)$, hopefully. (Note that by the last equation in (22), $p_{k+1} = p_k + F(z_{k+1})/\varepsilon_k$.)

To analyze the convergence of this iteration, we apply Theorem 2 to obtain an estimate of the form

$$\| p_{k+1} - p^* \|_{W^{1,p}} \leqq c \varepsilon_k \| p^h - p^* \|_{L^p} = c \varepsilon_k \| p_k - p^* \|_{L^p} \leqq c \varepsilon_k \| p_k - p^* \|_{W^{1,p}}.$$

If $c \varepsilon_k \leqq r < 1$ for every $k$, then

$$\| p_k - p^* \|_{W^{1,p}} \leqq r^k \| p_0 - p^* \|_{L^p}.$$

Moreover, by Theorem 2, we have

$$\| x_k - x^* \|_{W^{1,p}} + \| u_k - u^* \|_{L^\infty} \leqq r^k \| p_0 - p^* \|_{L^p}.$$

These observations are summarized in Corollary 1.

COROLLARY 1. *Under the hypotheses of Theorem 2 and for $\varepsilon_0 \| p_0 - p^* \|$ and $\sup \varepsilon_k$ sufficiently small, there exists a sequence* $(x_k, u_k, p_k)$, $k = 1, 2, \cdots$, *with the following properties:* $(x, u, p) = (x_k, u_k, p_k)$ *satisfies* (22) *when $\varepsilon = \varepsilon_{k-1}$ and $p^h = p_{k-1}$, and*

$$\| x_k - x^* \|_{W^{1,p}} + \| u_k - u^* \|_{L^\infty} + \| p_k - p^* \|_{W^{1,p}} \leqq r^k \| p_0 - p^* \|_{L^p}$$

*for some $0 \leqq r < 1$ and for $k = 1, 2, \cdots$.*

**5. Optimality.** We now show that for $\varepsilon$ sufficiently small, the $x_\varepsilon$ and $u_\varepsilon$ given by Theorem 2 locally minimize the augmented Lagrangian (2). We begin by stating a result whose proof is contained in Theorem 2.5 and Lemma 2.6 of [17].

LEMMA 5. *Let $\pi$ be a symmetric, continuous bilinear form defined on a Hilbert space $V$ with inner product $\langle \cdot, \cdot \rangle_V$, let $K$ be a convex subset of $V$, and let $L: V \to V$ be a linear map. Suppose that there exist positive $\rho$ and $\alpha$ such that*

$$(27) \qquad \pi(v, v) \geqq \alpha \langle v, v \rangle_V \quad \text{for all } v \in K \text{ with } L(v) = 0,$$

*and that any $v \in K$ can be expressed $v = v_1 + v_2$, where $v_1 \in K$, $L(v_1) = 0$, and $\| L v_2 \| \geqq \rho \| v_2 \|$. If $\alpha \rho^2 > \varepsilon \| \pi \| (\alpha + \| \pi \|)$, where*

$$\| \pi \| = \sup_{\substack{v \in V \\ \| v \| = 1}} \sup_{\substack{w \in V \\ \| w \| = 1}} \pi(v, w),$$

*then there exists $\beta > 0$ such that*

$$\pi(v, v) + \frac{1}{\varepsilon}\langle L(v), L(v)\rangle \geqq \beta\langle v, v\rangle_V \quad \text{for every } v \in K.$$

*The constant $\beta$ is independent of $\varepsilon$ for $\varepsilon$ sufficiently small.*

Assuming that the coercivity assumption (8) holds, we apply Lemma 5 to the linear-quadratic problem of § 2. If $z = (x, u) \in H^1 \times L^2$, $x(0) = 0$, and $u = v_1 - v_2$ for some $v_1$ and $v_2 \in U$, then we can write $z = z_1 + z_2$, where $z_1 = (x_1, u_1)$ with $u_1 = u$ and $x_1$ the solution to

$$M(x_1, u_1) = 0 \quad \text{and} \quad x_1(0) = 0,$$

and $z_2 = (x_2, u_2)$ with $u_2 = 0$ and $x_2 = x - x_1$. Since $u_2 = 0$, it follows that $M(z_2) = Ax_2 - \dot{x}_2$. Since $u_2 = 0$ and $x_2(0) = 0$, there exists $\rho > 0$ such that

$$(28) \qquad \langle M(z_2), M(z_2)\rangle \geqq \rho^2[\langle x_2, x_2\rangle_{H^1} + \langle u_2, u_2\rangle].$$

Lemma 5 and (8) and (28) give us the following lemma.

LEMMA 6. *If (8) holds, then for some positive constants $\bar{\varepsilon}$ and $\beta$, we have*

$$(29) \qquad \langle z, Pz\rangle + \frac{1}{\varepsilon}\langle M(z), M(z)\rangle \geqq \beta[\langle x, x\rangle_{H^1} + \langle u, u\rangle]$$

*whenever $0 < \varepsilon \leqq \bar{\varepsilon}$ and $z = (x, u) \in H^1 \times L^2$ with $x(0) = 0$ and $u = u_1 - u_2$ for some $u_1$ and $u_2 \in U$.*

Returning to the general augmented Lagrangian (2), we use Lemma 6 to prove local strict convexity of the cost functional.

LEMMA 7. *Suppose that (1) has a solution $z^* = (x^*, u^*)$, that there exists a bounded open set $\Sigma \subset R^{n+m}$, where both $f$ and $h$ are twice continuously differentiable, and that there exists a $\delta > 0$ with $z^*(t) \in \Sigma$ and the distance from $z^*(t)$ to the boundary of $\Sigma$ at least $\delta$ for almost every $t \in I$. If (8) holds, then for some positive $\gamma$ and $\bar{\varepsilon}$, and for every $z$ and $p$ in an $L^\infty$ neighborhood of $z^*$ and $p^*$, we have*

$$(30) \quad \langle \delta z, H_{zz}(z, p)\delta z\rangle + \frac{1}{\varepsilon}\langle f_z(z)\delta z - \dot{\delta x}, f_z(z)\delta z - \dot{\delta x}\rangle \geqq \gamma[\langle \delta x, \delta x\rangle_{H^1} + \langle \delta u, \delta u\rangle]$$

*whenever $0 < \varepsilon \leqq \bar{\varepsilon}$ and $\delta z = (\delta x, \delta u) \in H^1 \times L^2$ with $\delta x(0) = 0$ and $\delta u = u_1 - u_2$ for some $u_1$ and $u_2 \in U$.*

*Proof.* For $z = z^*$ and $p = p^*$, it follows from Lemma 6 that there exist positive constants $\beta$ and $\bar{\varepsilon}$ such that

$$(31) \quad \langle \delta z, H_{zz}(z^*, p^*)\delta z\rangle + \frac{1}{\varepsilon}\langle f_z(z^*)\delta z - \dot{\delta x}, f_z(z^*)\delta z - \dot{\delta x}\rangle \geqq \beta[\langle \delta x, \delta x\rangle_{H^1} + \langle \delta u, \delta u\rangle]$$

whenever $0 < \varepsilon \leqq \bar{\varepsilon}$ and $\delta z = (\delta x, \delta u) \in H^1 \times L^2$ with $\delta x(0) = 0$ and $\delta u = u_1 - u_2$ for some $u_1$ and $u_2 \in U$. Since decreasing $\varepsilon$ increases the second term on the left side of (30), let us examine the expression

$$\langle \delta z, H_{zz}(z, p)\delta z\rangle + \frac{1}{\bar{\varepsilon}}\langle f_z(z)\delta z - \dot{\delta x}, f_z(z)\delta z - \dot{\delta x}\rangle.$$

From (31) it follows that, for $z$ and $p$ in an $L^\infty$ neighborhood of $z^*$ and $p^*$, we have

$$\langle \delta z, H_{zz}(z, p)\delta z\rangle + \frac{1}{\bar{\varepsilon}}\langle f_z(z)\delta z - \dot{\delta x}, f_z(z)\delta z - \dot{\delta x}\rangle \geqq \frac{\beta}{2}[\langle \delta x, \delta x\rangle_{H^1} + \langle \delta u, \delta u\rangle].$$

Replacing $\bar{\varepsilon}$ by $\varepsilon$ completes the proof. $\quad\square$

THEOREM 3. *If the hypotheses of Theorem 2 hold, then for $\varepsilon\|p^h - p^*\|_{L^p}$ sufficiently small, the $x_\varepsilon$ and $u_\varepsilon$ given by Theorem 2 locally minimize the augmented Lagrangian (2).*

*Proof.* Let $N_z$ and $N_p$ denote $L^\infty$ neighborhoods of $z^*$ and $p^*$ where (30) holds. Choose $N_z$ small enough that it is contained in the sphere with center $z^*$ and radius $\delta$. Choose $\varepsilon \leqq \bar\varepsilon$ small enough that $z_\varepsilon \in N_z$ and $p_\varepsilon \in N_p$. Let $N_\varepsilon \subset N_z$ denote a neighborhood of $z_\varepsilon$ in $W^{1,\infty}$ for which $p^h + F(z)/\varepsilon$ lies in $N_p$ for each $z \in N_\varepsilon$. For every $z \in N_\varepsilon$ and $\delta z = (\delta x, \delta u) \in W^{1,\infty} \times L^\infty$, we have

$$
(32) \quad \frac{d^2}{d\tau^2} \left\{ C(z + \tau\delta z) + \langle p^h, F(z + \tau\delta z) \rangle + \frac{1}{2\varepsilon} \langle F(z + \tau\delta z), F(z + \tau\delta z) \rangle \right\}_{\tau=0}
$$

$$
= \langle \delta z, H_{zz}(z, p)\delta z \rangle + \frac{1}{\varepsilon} \langle f_z(z)\delta z - \dot{\delta x}, f_z(z)\delta z - \dot{\delta x} \rangle,
$$

where $p = p^h + F(z)/\varepsilon$. By Lemma 7 it follows that if $\delta x(0) = 0$ and $\delta u = u_1 - u_2$ for some $u_1$ and $u_2 \in U$, then the second derivative in (32) is positive whenever $\delta z \neq 0$. This positivity for the second derivative coupled with (19) for the first derivative implies that $z_\varepsilon$ is the unique minimizer of the augmented Lagrangian (2) over feasible points in $N_\varepsilon$.   □

**6. Finite-dimensional approximations.** In practice, the minimization of the augmented Lagrangian is carried out in a finite-dimensional space. In this section we determine the relationship between the dimension of the finite-dimensional space and the size of $\varepsilon$ so that the total error is minimized. To simplify the discussion, we drop the constraint $u(t) \in \Omega$ and set $p^h = 0$. That is, we consider the following augmented Lagrangian:

$$
(33) \quad \underset{z \in Z}{\text{minimize}} \; C(z) + \frac{1}{2\varepsilon} \langle F(z), F(z) \rangle,
$$

where $Z = \{(x, u) \in W^{1,\infty} \times L^\infty : x(0) = a\}$. Given a subspace $Z^h$ of $Z$, the approximating problem is

$$
(34) \quad \underset{z \in Z^h}{\text{minimize}} \; C(z) + \frac{1}{2\varepsilon} \langle F(z), F(z) \rangle.
$$

We will consider finite-element approximations in which case $h$ typically denotes the diameter of the largest mesh interval associated with the finite-element space.

First, let us observe that if $\varepsilon \to 0$ while $h$ is held fixed, the solution to (34) generally moves away from the solution to the original problem (1). The following simple problem illustrates this property:

$$
(35) \quad \begin{aligned} &\text{minimize } C(x, u) \\ &\text{subject to } \dot{x}(t) = x(t) + u(t) \quad \text{a.e. } t \in I, \quad x(0) = 1. \end{aligned}
$$

Partitioning the interval $I = [0, 1]$ into a uniform mesh, let us approximate $u$ by a piecewise constant function, and $x$ by a continuous, piecewise linear function. For fixed $h$, let $z_\varepsilon = (x_\varepsilon, u_\varepsilon)$ denote the minimizer in (34). As $\varepsilon \to 0$ in (34), the penalty term forces $F(z_\varepsilon) = x_\varepsilon + u_\varepsilon - \dot{x}_\varepsilon$ to zero. Since both $u_\varepsilon$ and the derivative of $x_\varepsilon$ are piecewise constant and

$$
x_\varepsilon = F(z_\varepsilon) + \dot{x}_\varepsilon - u_\varepsilon,
$$

it follows that as $\varepsilon \to 0$, $x_\varepsilon$ approaches a piecewise constant function. Since $x_\varepsilon$ is a continuous piecewise linear function and $x_\varepsilon(0) = 1$, we conclude that $x_\varepsilon$ approaches the function $x(t) = 1$ for every $t$. Moreover, if $F(z_\varepsilon) \to 0$ and $x_\varepsilon$ approaches 1, then $u_\varepsilon$ approaches the function $u(t) = -1$ for every $t \in [0, 1]$. On the other hand, $x(t) = 1$ and

$u(t) = -1$ for every $t \in [0, 1]$ is not the solution to (35) in general. Therefore, letting $\varepsilon \to 0$ while holding $h$ fixed moves us away from the desired solution.

The fundamental problem with letting $\varepsilon \to 0$ while holding $h$ fixed is that the null space of $M$, the linear system dynamics, is not "rich" enough to achieve a good approximation to $z^*$. That is, as $\varepsilon \to 0$, the augmented Lagrangian approximation to (35) approaches the null space of $M$ restricted to $Z^h$. However, for every choice of $h$, the null space of $M$ restricted to $Z^h$ is the single pair $(x, u)$, where $x(t) = 1$ and $u(t) = -1$ for every $t \in [0, 1]$. Since there is only one element in the null space of $M$ restricted to $Z^h$, we cannot achieve a good approximation to $z^*$ as $\varepsilon \to 0$. Chen et al. [6] give a necessary and sufficient condition under which the error in the solution to (34) can be decomposed into the sum of a term depending only on $h$ and a term depending only on $\varepsilon$. There is one exceptional case where it is possible to achieve good approximations even as $\varepsilon \to 0$: The system dynamics is $\dot{x} = u$ and the finite-element space used to approximate $u$ is the derivative of the finite-element space used to approximate $x$.

Nonetheless, as problem (35) indicates, $h$ generally needs to approach zero as $\varepsilon \to 0$ to ensure that the solution to (34) is a good approximation to the solution of the original problem. Let us now examine the optimal relation between $h$ and $\varepsilon$. We begin with a theoretical analysis for an abstract linear-quadratic problem:

(36)
$$\text{minimize } \tfrac{1}{2}\pi(v, v)$$
$$\text{subject to } L(v) = f,$$

where $\pi$ is a symmetric, continuous coercive bilinear form defined on a Hilbert space $V$, $\pi(v, v) \geqq \alpha \langle v, v \rangle_V$ for some $\alpha > 0$ and for every $v \in V$, and $L$ is a continuous linear operator whose range is a Hilbert space $W$ that contains $f$. The penalty approximation $v_\varepsilon$ to the solution $v^*$ of (36) is obtained by solving the problem

(37)
$$\underset{v \in V}{\text{minimize}} \; \pi(v, v) + \frac{1}{\varepsilon} \langle L(v) - f, L(v) - f \rangle_W.$$

The classical analysis of penalty approximations gives us an estimate of the form $\|v_\varepsilon - v^*\| \leqq c\varepsilon$.

Replacing the space $V$ of (37) by a finite-dimensional subspace $V^h$ yields the approximation

(38)
$$\underset{v \in V^h}{\text{minimize}} \; \pi(v, v) + \frac{1}{\varepsilon} \langle L(v) - f, L(v) - f \rangle_W.$$

If $v^h$ denotes the solution to (38), then the classical analysis of finite-element approximations (see [35]) gives the estimate

(39)
$$\pi_\varepsilon(v_\varepsilon - v^h, v_\varepsilon - v^h) \leqq \underset{v \in V^h}{\inf} \; \pi_\varepsilon(v_\varepsilon - v, v_\varepsilon - v),$$

where $\pi_\varepsilon(v, v) = \pi(v, v) + \varepsilon^{-1}\langle L(v), L(v) \rangle_W$. For standard finite-element spaces and regularity assumptions, it follows from (39) that

(40)
$$\alpha\|v_\varepsilon - v^h\|^2 + \varepsilon^{-1}\|L(v_\varepsilon - v^h)\|^2 \leqq \varepsilon^{-1}O(h^{2r}),$$

where $r$ is the degree of approximation associated with the finite-element space (when the control is approximated by a piecewise constant function and the state is approximated by a continuous, piecewise linear function, $r = 1$). Finally, by (40) and the triangle inequality, we conclude that

(41)
$$\|v^* - v^h\| \leqq \|v^* - v_\varepsilon\| + \|v_\varepsilon - v^h\| \leqq c\left[\varepsilon + \frac{h^r}{\sqrt{\varepsilon}}\right].$$

For fixed $h$, the right side of (41) is minimized by taking $\varepsilon = (h^r/2)^{2/3}$, and this choice for $\varepsilon$ gives $\|v^* - v^h\| = O(h^{2r/3})$.

In numerical experiments, it is observed that $\varepsilon = ch^{2r/3}$ is not optimal—the optimal $\varepsilon$ approaches zero even faster than $h^{2r/3}$ does, since the error term in (41) associated with the finite-element space overestimates the actual error. We now show rigorously that, for a linear quadratic problem with the state approximated by continuous piecewise linear finite elements and with the control approximated by piecewise constant finite elements, $\varepsilon = ch$ is optimal. For simplicity, the following problem is our model:

$$\text{(42)} \qquad \text{minimize} \quad \frac{1}{2} \int_I |x(t)|^2 + |u(t)|^2 \, dt$$

$$\text{subject to } M(z) = 0, \quad x(0) = a, \quad z = (x, u) \in H^1 \times L^2,$$

where $M(z) = Ax + Bu - \dot{x}$ with $A$ and $B$ constant, time-invariant matrices. The penalty approximation $z_\varepsilon$ to the solution $z^*$ of (42) is obtained by solving the problem

$$\text{(43)} \qquad \underset{\substack{z=(x,u)\in H^1 \times L^2 \\ x(0)=a}}{\text{minimize}} \ \pi_\varepsilon(z, z) \quad \text{where } \pi_\varepsilon(z, z) = \langle z, z \rangle + \frac{1}{\varepsilon} \langle M(z), M(z) \rangle.$$

Introducing a finite-dimensional subspace $Z^h = X^h \times U^h \subset H^1 \times L^2$, we have the following discrete approximation to (43):

$$\text{(44)} \qquad \underset{\substack{z=(x,u)\in Z^h \\ x(0)=a}}{\text{minimize}} \ \pi_\varepsilon(z, z).$$

Let $z^h = (x^h, u^h)$ denote the solution to (44) (which depends on $\varepsilon$) and let $z_\varepsilon = (x_\varepsilon, u_\varepsilon)$ denote the solution to (43). From the vanishing of the first variation, we have

$$\pi_\varepsilon(z_\varepsilon, \phi) = 0 \quad \text{for all } \phi \in H_0^1 \times L^2, \qquad \pi_\varepsilon(z^h, \phi) = 0 \quad \text{for all } \phi \in Z^h \cap H_0^1 \times L^2.$$

It follows that

$$\text{(45)} \qquad \pi_\varepsilon(z^h - z^I, \phi) = \pi_\varepsilon(z_\varepsilon - z^I, \phi) \quad \text{for all } z^I \in Z^h \text{ and } \phi \in Z^h \cap H_0^1 \times L^2.$$

Although the $x^I$ and $u^I$ components of $z^I$ are independent of each other, the distance between $x^h$ and $x^I$ can be bounded in terms of the distance between $u^h$ and $u^I$, and the distance between $z_\varepsilon$ and $z^I$, as stated in Lemma 8.

LEMMA 8. *If $z^h$ denotes the solution to (44) and $z_\varepsilon$ denotes the solution to (43), then for any $z^I = (x^I, u^I) \in Z^h$ with $x^I(0) = a$, we have*

$$\|M(z^h - z^I)\| \leq c\|z_\varepsilon - z^I\|_{H^1 \times L^2} \quad \text{and} \quad \|x^h - x^I\|_{H^1} \leq c\|z_\varepsilon - z^I\|_{H^1 \times L^2} + c\|u^h - u^I\|_{L^2}$$

*where $c$ is independent of $\varepsilon$ for $\varepsilon$ sufficiently small.*

*Proof.* Inserting $\phi = z^h - z^I$ in (45) gives

$$\text{(46)} \quad \varepsilon\|z^h - z^I\|^2 + \|M(z^h - z^I)\|^2 \leq \varepsilon\langle z_\varepsilon - z^I, z^h - z^I \rangle + \langle M(z_\varepsilon - z^I), M(z^h - z^I) \rangle.$$

From the relation $2\langle f, g \rangle \leq \langle f, f \rangle + \langle g, g \rangle$, we conclude that

$$\text{(47)} \qquad \|M(z^h - z^I)\|^2 \leq \varepsilon\|z_\varepsilon - z^I\|^2 + \|M(z_\varepsilon - z^I)\|^2,$$

which yields the first inequality of Lemma 8. To obtain the second inequality, we start with the relation

$$\|x\|_{H^1} \leq c\|\dot{x} - Ax\|_{L^2} = c\|Bu - M(x, u)\|_{L^2} \quad \text{for every } x \in H_0^1 \text{ and } u \in L^2.$$

Inserting $(x, u) = z^h - z^I$ gives

$$\|x^h - x^I\|_{H^1} \leq \|M(z^h - z^I)\| + c\|u^h - u^I\|.$$

The first inequality of Lemma 8 completes the proof.    □

In order to show that $\varepsilon = ch$ is optimal when the state is approximated by continuous, piecewise linear finite elements and when the control is approximated by piecewise constant finite elements, we assume that the meshes associated with $X^h$ and with $U^h$ are identical, quasi-uniform meshes. That is, the ratio between the width of the largest and the smallest mesh interval is bounded by a uniform constant. If $h$ denotes the width of the largest mesh interval, then the interpolation error has the following bounds (see [35]):

$$(48) \qquad \|x^I - x^*\|_{H^1} \leqq ch\|x^*\|_{H^2} \quad \text{and} \quad \|u^I - u^*\|_{L^2} \leqq ch\|u^*\|_{H^1}.$$

On the other hand, except for special choices of $x^*$ and $u^*$, a lower bound of the form

$$\underset{x^h \in X^h, u^h \in U^h}{\text{minimum}} \|x^h - x^*\|_{H^1} + \|u^h - u^*\|_{L^2} \geqq ch,$$

$c > 0$, also holds. Consequently, if we can show that

$$(49) \qquad \|x^h - x^*\|_{H^1} + \|u^h - u^*\|_{L^2} = O(h)$$

for $\varepsilon = ch$, then $\varepsilon = ch$ is optimal in the sense that the error is (asymptotically) as small as possible.

THEOREM 4. *If $X^h$ consists of continuous piecewise linear polynomials and $U^h$ consists of piecewise constant polynomials, then for a quasi-uniform mesh and $\varepsilon = ch$ with $c > 0$, estimate (49) holds.*

Proof. By Theorem 2,

$$\|x_\varepsilon - x^*\|_{H^1} + \|u_\varepsilon - u^*\|_{L^2} + \|p_\varepsilon - p^*\|_{H^1} = O(\varepsilon) = O(h)$$

since $\varepsilon = ch$. From the necessary condition (22) associated with (43), we conclude that $p_\varepsilon, x_\varepsilon$, and $u_\varepsilon$ are infinitely differentiable with the norm of each derivative bounded uniformly in $\varepsilon$ for $\varepsilon$ sufficiently small. The distance between $z_\varepsilon$ and $z^h$ is estimated using (46) with $z^I$ constructed in the following way: $z^I = (x^I, u^I)$, where $u^I$ is the $L^2$ projection of $u_\varepsilon$ into $U^h$ and $x^I = x_1^I + \varepsilon x_2^I$, where $x_2^I$ interpolates the solution $x_2$ to

$$(50) \qquad \dot{x}_2 = Ax_2 - p_\varepsilon, \qquad x_2(0) = 0,$$

and $x = x_1^I$ is the solution to the problem

$$(51) \qquad \begin{aligned} &\underset{x \in X^h}{\text{minimize}} \|\dot{x}_1 - \dot{x}\|_{L^2} \\ &\text{subject to } \dot{x}_1 = Ax_1 + Bu_\varepsilon, \quad x_1(0) = a = x(0). \end{aligned}$$

By the interpolation error estimate (48) and by the classical theory (see [35]) for the error in elliptic projections, we have

$$(52) \qquad \|u_\varepsilon - u^I\|_{L^2} = O(h) = \|x_1 - x_1^I\|_{H^1} = O(h) = \|x_2 - x_2^I\|_{H^1}.$$

Since $x_\varepsilon$ satisfies (22) while $x_1$ and $x_2$ satisfy the differential equations in (50) and (51), it follows that $x_\varepsilon = x_1 + \varepsilon x_2$ and

$$(53) \qquad \|x_\varepsilon - x^I\|_{H^1} \leqq \|x_1 - x_1^I\|_{H^1} + \varepsilon\|x_2 - x_2^I\|_{H^1} = O(h).$$

Below, we will derive the following bound for the last term in (46):

$$(54) \qquad \langle M(z_\varepsilon - z^I), M(z^h - z^I) \rangle \leqq O(h^3) + O(h^2)\|u^h - u^I\|.$$

Consequently, by (46) and the inequality $2\langle f, g \rangle \leqq \langle f, f \rangle + \langle g, g \rangle$, we have

$$\varepsilon\|z^h - z^I\|^2 + 2\|M(z^h - z^I)\|^2 \leqq \varepsilon\|z_\varepsilon - z^I\|^2 + O(h^3) + O(h^2)\|u^h - u^I\|.$$

Dividing by $\varepsilon$ and putting $\varepsilon = ch$ leads to

$$\|u^h - u^I\|^2 \leqq \|z_\varepsilon - z^I\|^2 + O(h^2) + O(h)\|u^h - u^I\|,$$

which yields $\|u^h - u^I\| = O(h)$ from (52) and (53). Moreover, by Lemma 8, we have

$$\|x^h - x^I\|_{H^1} \leqq c\|z_\varepsilon - z^I\|_{H^1 \times L^2} + c\|u^h - u^I\|_{L^2} = O(h).$$

Finally, the triangle inequality and Theorem 2 complete the proof:

$$\|x^h - x^*\|_{H^1} \leqq \|x^h - x^I\|_{H^1} + \|x^I - x_\varepsilon\|_{H^1} + \|x_\varepsilon - x^*\|_{H^1} = O(h),$$

$$\|u^h - u^*\|_{L^2} \leqq \|u^h - u^I\|_{L^2} + \|u^I - u_\varepsilon\|_{L^2} + \|u_\varepsilon - u^*\|_{L^2} = O(h).$$

Let us now verify (54). Since $Mz_\varepsilon = \varepsilon p_\varepsilon = \varepsilon(Ax_2 - \dot{x}_2)$ and since

$$M(x_1^I, u^I) = Bu^I + Ax_1^I - \dot{x}_1^I = B(u^I - u_\varepsilon) + A(x_1^I - x_1) - (\dot{x}_1^I - \dot{x}_1),$$

it follows that

(55) $\quad M(z_\varepsilon - z^I) = \varepsilon[(\dot{x}_2^I - \dot{x}_2) - A(x_2^I - x_2)] + [(\dot{x}_1^I - \dot{x}_1) - A(x_1^I - x_1) - B(u^I - u_\varepsilon)].$

Hence, the left side of (54) decomposes into three terms:

$$\langle M(z_\varepsilon - z^I), M(z^h - z^I) \rangle = T_1 + T_2 - T_3,$$

where

$$T_1 = \langle \varepsilon[(\dot{x}_2^I - \dot{x}_2) - A(x_2^I - x_2)] - A(x_1^I - x_1), M(z^h - z^I) \rangle,$$

$$T_2 = \langle (\dot{x}_1^I - \dot{x}_1) - B(u^I - u_\varepsilon), (\dot{x}^I - \dot{x}^h) - B(u^I - u^h) \rangle,$$

$$T_3 = \langle (\dot{x}_1^I - \dot{x}_1) - B(u^I - u_\varepsilon), A(x^I - x^h) \rangle.$$

Starting with the first term, we apply the interpolation error bound (52) to obtain

(56) $$\|(\dot{x}_2^I - \dot{x}_2) - A(x_2^I - x_2)\| = O(h).$$

By the Aubin–Nitsche duality trick (see [35]) for estimating the $L^2$ error in an elliptic projection,

(57) $$\|A(x_1^I - x_1)\| = O(h^2).$$

Combining (56), (57), and the first inequality of Lemma 8, we conclude that $T_1 = O(h^3)$. Since $u^I$ is the orthogonal projection of $u_\varepsilon$ into $U^h$, $u^I - u_\varepsilon$ is orthogonal to the space of piecewise constant functions. Similarly, from the structure of the minimization problem used to generate $x_1^I$, the difference $\dot{x}_1 - \dot{x}_1^I$ is orthogonal to all piecewise constant functions. Thus $T_2 = 0$. Finally, let $q$ denote the projection of $x^I - x^h$ into the space of piecewise constant functions. Again, exploiting orthogonality, we have

$$\langle (\dot{x}_1^I - \dot{x}_1) - B(u^I - u_\varepsilon), Aq \rangle = 0,$$

so that $T_3$ can be expressed

$$T_3 = \langle (\dot{x}_1^I - \dot{x}_1) - B(u^I - u_\varepsilon), A(x^I - x^h - q) \rangle.$$

By the Schwarz inequality, the interpolation error estimate (48), and Lemma 8, we have

$$T_3 \leqq ch^2\|x^I - x^h\|_{H_1} \leqq ch^3 + ch^2\|u^h - u^I\|.$$

Combining these bounds for $T_1$, $T_2$ and $T_3$ yields (54). $\quad\square$

*Remark.* Although the size of the penalty parameter is crucial when an augmented Lagrangian is discretized, it is less crucial if an optimal control problem is discretized and the discrete problem is solved by a multiplier method. As $\varepsilon \to 0$ in the augmented

Lagrangian for the discrete problem, the solution to the penalized problem typically converges to a solution of the discretized problem. And for an appropriate discretization (for example, see [8], [9], or [14] and the references they cite), the solution to the discrete problem converges to the solution of the continuous problem as the mesh is refined. Hence, letting $\varepsilon \to 0$ while fixing the discretization does not interfere with the convergence when the discretization is "appropriate."

**7. Numerical experiments.** The dependence between the optimal $\varepsilon$ and $h$ was investigated using four problems. The first problem was

(P1)

$$\text{minimize} \quad \frac{1}{2}\left\{x(1)^2 + \int_0^1 u(t)^2 \, dt\right\}$$

$$\text{subject to} \quad \dot{x}(t) = x(t) + u(t), \quad x(0) = 1,$$

with the optimal solution

$$x^*(t) = e^t + \frac{b}{2}[e^{-t} - e^t], \quad u^*(t) = -be^{-t}, \quad b = \frac{2e^2}{1+e^2}.$$

The second problem was

(P2)

$$\text{minimize} \quad \frac{1}{4}\int_0^1 2x(t)^2 + u(t)^2 \, dt$$

$$\text{subject to} \quad \dot{x}(t) = .5x(t) + u(t), \quad x(0) = 1,$$

with the optimal solution

$$x^*(t) = \frac{c\, e^{1.5t} - e^{-1.5t}}{c-1}, \quad u^*(t) = \frac{2\, e^{-1.5t} + c\, e^{1.5t}}{c-1}, \quad c = -2\, e^{-3}.$$

The third problem was

(P3)

$$\text{minimize} \quad \frac{1}{4}\int_0^1 1.25x(t)^2 + x(t)u(t) + u(t)^2 \, dt$$

$$\text{subject to} \quad \dot{x}(t) = .5x(t) + u(t), \quad x(0) = 1,$$

with the optimal solution

$$x^*(t) = \frac{\cosh(1-t)}{\cosh(1)}, \quad u^*(t) = -x^*(t)[\tanh(1-t) + .5].$$

The fourth problem was

(P4)

$$\text{minimize} \quad \frac{1}{2}\int_0^1 e^{-2t}x(t)^2 + u(t)^2 \, dt$$

$$\text{subject to} \quad \dot{x}(t) = x(t) + e^t u(t), \quad x(0) = \frac{1+3e}{2(1-e)}, \quad u(t) \le 1,$$

with the optimal solution

$$x^*(t) = e^t(t + x(0)), \quad u^*(t) = 1, \quad 0 \le t \le \tfrac{1}{2},$$

$$x^*(t) = (e^{2t} - e^2)/d, \quad u^*(t) = (e^t - e^{2-t})/d, \quad \tfrac{1}{2} \le t \le 1,$$

where $d = \sqrt{e}(1-e)$. Note that (P1)–(P3) are unconstrained quadratic cost problems while (P4) has a control constraint.

TABLE 1
The optimal $\varepsilon$ and error for the test problems.

| | P1 | | P2 | | P3 | | P4 | |
|---|---|---|---|---|---|---|---|---|
| $1/h$ | $1/\varepsilon$ | Error | $1/\varepsilon$ | Error | $1/\varepsilon$ | Error | $1/\varepsilon$ | Error |
| 10 | 54.4 | .02442 | 41.3 | .00959 | 32.9 | .00522 | 9.6 | .31600 |
| 20 | 108.5 | .01265 | 81.5 | .00525 | 66.1 | .00278 | 17.6 | .16840 |
| 40 | 216.7 | .00644 | 161.8 | .00275 | 132.6 | .00144 | 34.0 | .08681 |
| 80 | 433.1 | .00325 | 322.5 | .00141 | 265.4 | .00073 | 66.8 | .04405 |
| 160 | 866.1 | .00163 | 643.2 | .00071 | 530.4 | .00037 | 132.4 | .02219 |
| 320 | 1731.7 | .00082 | 1286.4 | .00036 | 1065.0 | .00018 | 263.5 | .01113 |

For the finite-dimensional problem (34), we employed a uniform mesh with mesh spacing $h$. The controls were approximated by piecewise constant polynomials while the states were approximated by continuous, piecewise linear polynomials. Table 1 gives the optimal $\varepsilon$ corresponding to various choices of $h$ (actually $1/\varepsilon$ is given for various choices of $1/h$). The optimal $\varepsilon$ was chosen to minimize the expression

$$(58) \qquad \langle x^* - x^h, x^* - x^h \rangle_{H^1} + \langle u^* - u^h, u^* - u^h \rangle,$$

where $(x^h, u^h)$ denotes the solution to (34) (which depends on $\varepsilon$). The column labeled "Error" in Table 1 gives the square root of the optimal expression (58). Clearly, both the optimal $\varepsilon$ and the optimal error are asymptotically proportional to $h$.

REFERENCES

[1] K. J. ARROW AND R. M. SOLOW, Gradient methods for constrained maxima, with weakened assumptions, in Studies in Linear and Nonlinear Programming, K. Arrow, L. Hurwicz, and H. Uzawa, eds., Stanford University Press, Stanford, CA, 1958.
[2] I. BABUŠKA, The finite element method with penalty, Math. Comp., 27 (1973), pp. 221-228.
[3] A. V. BALAKRISHNAN, On a new computing technique in optimal control, SIAM J. Control, 6 (1968), pp. 149-173.
[4] D. P. BERTSEKAS, Constrained optimization and Lagrange multiplier methods, Academic Press, New York, 1982.
[5] G. CHEN AND W. H. MILLS, Finite elements and terminal penalization for quadratic cost optimal control problems governed by ordinary differential equations, SIAM J. Control Optim., 19 (1981), pp. 744-764.
[6] G. CHEN, W. H. MILLS, JR., S. SUN, AND D. A. YOST, Sharp error estimates for a finite element-penalty approach to a class of regulator problems, Math. Comp., 40 (1983), pp. 151-173.
[7] R. COURANT, Variational methods for the solution of problems of equilibrium and vibrations, Bull. Amer. Math. Soc., 49 (1943), pp. 1-23.
[8] J. CULLUM, Discrete approximations to continuous optimal control problems, SIAM J. Control, 7 (1969), pp. 32-49.
[9] ———, An explicit procedure for discretizing continuous, optimal control problems, J. Optim. Theory. Appl., 8 (1971), pp. 15-34.
[10] R. S. FALK, An analysis of the penalty method and extrapolation for the stationary Stokes equations, in Advances in Computer Methods for Partial Differential Equations, R. Vichnevetsky, ed., AICA, New Brunswick, NJ, 1975, pp. 66-69.
[11] R. S. FALK AND J. T. KING, A penalty and extrapolation method for the stationary Stokes equations, SIAM J. Numer. Anal., 13 (1976), pp. 814-829.

[12]  R. FORTIN AND R. GLOWINSKI, *Augmented Lagrangian methods: applications to the numerical solution of boundary-value problems*, North-Holland, Amsterdam, New York, 1983.

[13]  W. W. HAGER, *The Ritz-Trefftz method for state and control constrained optimal control problems*, SIAM J. Numer. Anal., 12 (1975), pp. 854-867.

[14]  ———, *Rates of convergence for discrete approximations to unconstrained control problems*, SIAM J. Numer. Anal., 13 (1976), pp. 449-472.

[15a]  ———, *Convex control and dual approximations. Part I*, Control Cybernet., 8 (1979), pp. 5-22.

[15b]  ———, *Convex control and dual approximations. Part II*, Control Cybernet., 8 (1979), pp. 73-86.

[16]  ———, *Dual approximations in optimal control*, SIAM J. Control Optim., 22 (1984), pp. 423-465.

[17]  ———, *Approximations to the multiplier method*, SIAM J. Numer. Anal., 22 (1985), pp. 16-46.

[18]  W. W. HAGER AND G. D. IANCULESCU, *Semi-dual approximations in optimal control: quadratic cost*, in Free Boundary Problems, Vol. II, Pavia, 1979, Ist. Naz. Alta Mat. Francesco Severi, Rome, 1980, pp. 321-332.

[19]  M. R. HESTENES, *Multiplier and gradient methods*, J. Optim. Theory Appl., 4 (1969), pp. 303-320.

[20]  E. ISAACSON AND H. B. KELLER, *Analysis of numerical methods*, John Wiley, New York, 1966.

[21]  H. KHESHGI AND M. LUSKIN, *On the variable sign penalty approximation of the Navier-Stokes equation*, in Nonlinear Partial Differential Equations, J. A. Smoller, ed., American Mathematical Society, Providence, RI, 1983, pp. 91-108.

[22]  N. KIKUCHI, *Convergence of a penalty-finite element approximation for an obstacle problem*, Numer. Math., 37 (1981), pp. 105-120.

[23]  J. T. KING, *New error bounds for the penalty method and extrapolation*, Numer. Math., 23 (1974), pp. 153-165.

[24]  J. L. LIONS, *Optimal control of systems governed by partial differential equations*, S. K. Mitter, trans., Springer-Verlag, Berlin, New York, 1971.

[25]  B. T. POLYAK, *The convergence rate of the penalty function method*, Zh. Vychisl. Mat. i Mat. Fiz., 11 (1971), pp. 3-11. (In Russian.) USSR Comput. Math. and Math. Phys., 11 (1971), pp. 1-12. (In English.)

[26]  B. T. POLYAK AND N. V. TRET'YAKOV, *The method of penalty estimates for conditional extremum problems*, Zh. Vychisl. Mat. i Mat. Fiz., 13 (1973), pp. 34-46. (In Russian.) USSR Comput. Math. and Math. Phys., 13 (1973), pp. 42-58. (In English.)

[27]  L. S. PONTRYAGIN, V. G. BOLTYANSKII, R. V. GAMKRELIDZE, AND E. F. MISHCHENKO, *The Mathematical Theory of Optimal Processes*, Wiley-Interscience, New York, 1962.

[28]  M. J. D POWELL, *A method for nonlinear constraints in minimization problems*, in Optimization, R. Fletcher, ed., Academic Press, New York, 1972.

[29]  J. N. REDDY, *On penalty function methods in the finite-element analysis of flow problems*, Internat. J. Numer. Methods Fluids, 2 (1982), pp. 151-171.

[30]  S. M. ROBINSON, *Strongly regular generalized equations*, Math. Oper. Res., 5 (1980), pp. 43-62.

[31]  R. T. ROCKAFELLAR, *The multiplier method of Hestenes and Powell applied to convex programming*, J. Optim. Theory Appl., 12 (1973), pp. 555-562.

[32]  ———, *A dual approach to solving nonlinear programming problems by unconstrained optimization*, Math. Programming, 5 (1973), pp. 354-373.

[33]  ———, *Augmented Lagrange multiplier functions and duality in nonconvex programming*, SIAM J. Control, 12 (1974), pp. 268-285.

[34]  D. L. RUSSELL, *Penalty functions and bounded phase coordinate control*, SIAM J. Control, 2 (1965), pp. 409-422.

[35]  G. STRANG AND G. J. FIX, *An analysis of the finite element method*, Prentice-Hall, Englewood Cliffs, NJ, 1973.

[36]  W. WERNER, *Penalty function methods for the numerical solution of nonlinear obstacle problems with finite elements*, Z. Angew. Math. Mech., 61 (1981), pp. 133-139.