

UNASA-CR-172,301

**NASA Contractor Report** 172301

# ICASE

**FOR REFERENCE**

**NOT TO BE TAKEN FROM THIS ROOM**

ESTIMATION OF DISCONTINUOUS COEFFICIENTS IN PARABOLIC SYSTEMS:  
APPLICATIONS TO RESERVOIR SIMULATION

Patricia Daniel Lamm

NASA-CR-172301  
19840012205

Contract Nos. NAS1-17130, NAS1-16394  
February 1984

INSTITUTE FOR COMPUTER APPLICATIONS IN SCIENCE AND ENGINEERING  
NASA Langley Research Center, Hampton, Virginia 23665

Operated by the Universities Space Research Association

**NASA**

National Aeronautics and  
Space Administration

**Langley Research Center**  
Hampton, Virginia 23665

**LIBRARY COPY**

FEB 20 1984

LANGLEY RESEARCH CENTER  
LIBRARY, NASA  
HAMPTON, VIRGINIA



8 1 1 RN/NASA-CR-172301

DISPLAY 08/2/1

84N20273\*\* ISSUE 10 PAGE 1568 CATEGORY 64 RPT#: NASA-CR-172301 NAS  
1.26:172301 CNT#: NAS1-17130 NAS1-16394 MAG1-258 NSF MCS-82-00883  
84/02/00 56 PAGES UNCLASSIFIED DOCUMENT

UTTL: Estimation of discontinuous coefficients in parabolic systems:  
Applications to reservoir simulation TLSP: Final Report

AUTH: A/LAMM, P. D.

CORP: Southern Methodist Univ., Dallas, Tex. AVAIL. NTIS SAP: HC A04/MF A01

MAJS: /\*COEFFICIENTS/\*DISCONTINUITY/\*PARABOLAS

MINS: / ALGORITHMS/ CONVERGENCE/ SPLINE FUNCTIONS

ABA: Author

ABS: Spline based techniques for estimating spatially varying parameters that  
appear in parabolic distributed systems (typical of those found in  
reservoir simulation problems) are presented. The problem of determining  
discontinuous coefficients, estimating both the functional shape and  
points of discontinuity for such parameters is discussed. Convergence  
results and a summary of numerical performance of the resulting algorithms  
are given.

ENTER:



# ESTIMATION OF DISCONTINUOUS COEFFICIENTS IN PARABOLIC SYSTEMS:

## APPLICATIONS TO RESERVOIR SIMULATION

Patricia Daniel Lamm  
Southern Methodist University

### ABSTRACT

We present spline-based techniques for estimating spatially varying parameters that appear in parabolic distributed systems (typical of those found in reservoir simulation problems). In particular, we discuss the problem of determining discontinuous coefficients, estimating both the functional shape and points of discontinuity for such parameters. In addition, our ideas may also be applied to problems with unknown initial conditions and unknown parameters appearing in terms representing external forces. Convergence results and a summary of numerical performance of the resulting algorithms are given.

---

Research supported in part by NSF Grant MCS-8200883 and NASA Grant NAG-1-258. Part of the research was carried out while the author was a visitor at the Institute for Computer Applications in Science and Engineering (ICASE), NASA Langley Research Center, Hampton, VA., which is operated under NASA Contract Nos. NAS1-17130 and NAS1-16394.



## 1. Introduction

We present here our efforts related to the estimation of discontinuous spatially varying coefficients in parabolic distributed systems. Although our ideas are applicable to a wide class of problems in which the determination of discontinuous coefficients is of importance (e.g., the propagation of waves through layered media; the dynamics of beams with "discontinuous" elastic properties), our work here is motivated by an inverse problem in reservoir simulation commonly referred to as "history matching". The problem in this case is to determine unknown parameters (such as permeability, porosity) that appear as coefficients in model reservoir equations. "Optimal" choices of these parameters should provide the best match between the observed and simulated production history at one or more wells. Information about these coefficients (functional shape and location of discontinuities) provides insight into physical properties of the reservoir and can indicate the location of abrupt structural changes; in addition, precise determination of these parameters is essential to the process of accurately simulating and predicting reservoir behavior.

The governing reservoir equations describe mathematically the physical and chemical processes occurring during primary hydrocarbon recovery or during enhanced recovery efforts (secondary or tertiary forms of recovery). Mathematical models vary widely depending on the physical process being described (miscible or immiscible fluid flow, thermal or fluid injection, etc.) and the types of observations available. Common to each model however is a system of rate equations (derived from Darcy's law, which relates flow rate to fluid pressure gradients) as well as appropriate conservation laws and equations of state. The resulting dynamical system is typically distributed in nature and of parabolic type [17], [18]; unknown parameters quite often include the porosity of

surrounding rock, or the ratio of pore volume to total volume, and (relative) permeability, which is the ability of the rock to transmit fluid [18]. Due to spatial changes in underground structure, it is highly likely that these parameters will vary spatially and contain numerous discontinuities.

In order to solve the inverse problem, data in the form of fluid pressure (or flow rate) is collected at the wells and used in a numerical parameter estimation process. There have been a large number of substantial contributors to the development of theoretical concepts and numerical algorithms for the history matching problem. An exhaustive list of related references would be too lengthy to include here; instead we refer the reader to [18] for an excellent survey of the outstanding efforts in this area. One numerical approach commonly taken involves subdividing the reservoir into a grid of smaller blocks; constant-valued parameters (which are allowed to vary independently from block to block) are then estimated. Unfortunately, if accurate solutions are desired, the grid size often must be quite small and thus the number of unknown parameters, as well as the dimension of the state space, can be very large--as many as 50,000 parameters or more [17]. (This is an unfortunate consequence of the fact that the parameters of interest--as well as state variables--are infinite-dimensional yet computations must be performed in a finite-dimensional setting.) Our goal here is to avoid some of the difficulties associated with the approach described above. Specifically, our ideas involve separating the order of state approximation from that of parameter estimation, so that the need for an approximate state space of high dimension does not impose the same requirements on the dimension of an approximate parameter space; this is accomplished by searching for parameters in classes of functions with quite general spatially varying representations. In order to focus attention on the problems associated



with estimating spatially varying discontinuous coefficients in this context, we consider an archetypical model of (parabolic) distributed type that admittedly is a simplified version of the fluid pressure equations associated with reservoir simulation (see [17], [18], and the references therein); nevertheless the model selected here is a prototype that contains the essential parameter-dependent terms for which we may begin our investigations. In the sections that follow we define the model equations of interest and construct an approximation framework in which we wish to consider the parameter estimation problem. Convergence results are presented for problems associated with either spatially distributed or "discrete" sample data. Finally, we discuss numerical implementation in general, and in the context of particular examples.

It is our intent in this report to examine convergence properties and implementation problems associated with these methods; we do not address such important questions as identifiability, observability, or general underlying properties of the governing partial differential equation system.

The notation used throughout is standard: For  $I \subseteq \mathbb{R}$  (the real line), we shall denote by  $C(I; X)$  the space of continuous functions  $f: I \rightarrow X$  with uniform norm  $|\cdot|_\infty$ ; by  $L_2(I; X)$  we mean the usual space of square-integrable "functions"  $f: I \rightarrow X$  with  $L_2$  norm  $|\cdot|_{L_2(I; X)}$  and inner product  $\langle \cdot, \cdot \rangle_{L_2(I; X)}$ . The Sobolev spaces  $H^p(I; X)$  and  $H_0^p(I; X)$  are defined as usual (see, for example, [1]). Whenever  $X = \mathbb{R}$ , we shall simplify notation by writing  $C(I)$  and  $L_2(I)$ , respectively, and, where no confusion results, by writing  $|\cdot|$  (and  $\langle \cdot, \cdot \rangle$ ) for the norm (and inner product) on  $L_2(0, 1)$ . In addition, no notational distinction will be made between a function  $f: I \rightarrow \mathbb{R}$  and its restriction to  $I_1 \subseteq I$ .

## 2. The Parameter Estimation Problem

As our fundamental state system we consider the scalar parabolic distributed system

$$(2.1) \quad \begin{cases} \frac{\partial u}{\partial t}(t, x) = \frac{1}{\rho(x)} \frac{\partial}{\partial x} (q(x) \frac{\partial u}{\partial x}(t, x)) + f(t, x; r(x)), & (t, x) \in (0, T) \times (0, 1), \\ u(t, 0) = u(t, 1) = 0, \\ u(0, x) = u_0(x). \end{cases}$$

Here  $q$  and  $\rho$  are discontinuous (positive) functions representing the permeability and porosity properties, respectively, of the fluid and surrounding rock; the points of discontinuity in these functions correspond to abrupt spatial changes in the physical flow region (such as might be associated with layered media). Both  $q$  and  $\rho$  are typically unknown so we shall consider the problem of estimating these parameters, as well as the function  $r$ ,  $r(x) \in R^P$  and the initial condition  $u_0$ , from observations of the state variable  $u$ .

For ease of presentation in the arguments that follow, it is assumed that  $\rho \equiv 1$ , although it is not difficult to extend our ideas to the case of non-constant (and unknown)  $\rho$ . We detail in Remark 3.1 the minor modifications one must make in the calculations found below in order to treat  $\rho$  as a functional parameter throughout.

To simplify notation, we assume that  $q$  is discontinuous at one point only,  $x = \xi$ , and that  $q$  is represented by

$$q = \phi_1 + H_\xi \phi_2$$

where  $\phi_1$  and  $\phi_2$  are continuous on  $[0,1]$ ; here  $H_\xi$  is the usual Heaviside function on  $[0,1]$  given by  $H_\xi = 1$  on  $[\xi,1]$ ,  $H_\xi = 0$  otherwise. There is a straightforward extension of our ideas to the case where

$$q = \phi_1 + \sum_{i=2}^{\mu} H_{\xi_{i-1}} \phi_i$$

$0 = \xi_0 < \xi_1 < \xi_2 < \dots < \xi_\mu = 1$ , except that notational difficulties become excessive. (We later demonstrate our approximation and estimation techniques for multiple discontinuity problems in the section on numerical findings.) Given the parameterization chosen for  $q$  we define the parameter vector  $\gamma = (\xi, \phi_1, \phi_2, r, u_0) = (s, u_0)$  as an element of the parameter set  $\Gamma \subseteq \tilde{\Gamma} = \mathcal{S} \times L_2(0,1)$ , where, for  $m, \bar{m}$  fixed,

$$\mathcal{S}(m, \bar{m}) \equiv \{s = (\xi, \phi_1, \phi_2, r) \in \mathbb{R} \times C[0,1] \times C[0,1] \times L_2((0,1); \mathbb{R}^P) \mid \xi \in (0,1),$$

$$\phi_i \in C^1[0,1], \text{ and } 0 < m \leq \phi_i(x) \leq \bar{m}$$

$$\text{for } i = 1, 2, \text{ and } x \in [0,1]\}.$$

Concerning  $\Gamma$  and the applied force  $f$ , we make the following (standing) hypotheses:

(H1) The parameter set  $\Gamma$  is compact;

(H2) For every  $r \in L_2((0,1); \mathbb{R}^P)$ , the map  $t \mapsto f(t, \cdot; r(\cdot)) : [0, T] \rightarrow L_2(0,1)$  is Hölder continuous with exponent  $\alpha$ ,  $0 < \alpha < 1$ .

(H3) The map  $r \mapsto f(\cdot, \cdot; r(\cdot))$  is continuous from  $L_2((0,1); \mathbb{R}^P)$  to  $L_2((0,T) \times (0,1))$ .

The parameter estimation problem associated with (2.1) consists of finding a parameter  $\gamma^* \in \Gamma$  that is "optimal" in the sense of providing the best match between observed data and model solutions to (2.1). Although a number of criteria may be used to measure "fit to data," we consider first a least squares criterion  $J$  that is defined in conjunction with distributed data: That is, given distributed observations  $\hat{u}_i \in L_2(0,1)$  at discrete times  $t_i \in (0,T)$ ,  $i = 1, \dots, n$ , we seek  $\gamma^* \in \Gamma$  that minimizes

$$(2.2) \quad J(\gamma) = \sum_{i=1}^n \int_0^1 |\mathfrak{C}(t_i, x; \gamma)u(t_i, x; \gamma) - \hat{u}_i(x)|^2 dx$$

over all  $\gamma \in \Gamma$ . For each  $(t_i, x)$ , the output map  $\mathfrak{C}(t_i, x; \gamma) : R \rightarrow R$  is assumed to be continuous in  $\gamma$  and such that the mapping  $x \rightarrow \mathfrak{C}(t_i, x; \gamma)\psi(x)$  is in  $L_2(0,1)$  whenever  $\psi \in L_2(0,1)$ . We note that data generally is not available in the distributed form given here; often this difficulty can be handled by fitting a curve (using linear interpolation, for example) to discrete data.

We also treat the problem of truly discrete data, i.e.,  $\hat{u}_{ij} \in R$  is observed sample data at  $(t_i, x_j)$ ,  $j = 1, \dots, \tilde{n}$ . In this case the parameter estimation problem consists of determining  $\tilde{\gamma}^* \in \Gamma$  that minimizes a "pointwise" fit-to-data criterion,

$$(2.3) \quad \tilde{J}(\gamma) = \sum_{i=1}^n \sum_{j=1}^{\tilde{n}} |\mathfrak{C}(t_i, x_j; \gamma)u(t_i, x_j; \gamma) - \hat{u}_{ij}|^2$$

over  $\gamma \in \Gamma$ . The use of discrete sample data leads to increased technical detail and additional smoothness hypotheses on  $u_0$  and  $f$ . We defer consideration of this particular estimation problem until we have fully developed an approximation theory for distributed estimation problems (i.e., identification problems with cost functional  $J$  defined in (2.2)); our findings for the "pointwise" estimation problem (involving  $\tilde{J}$ ) are then summarized in section 3.1 below.

Before we consider either of these estimation problems (where  $\gamma \in \Gamma$  is unknown and to be determined), we shall first consider the existence of solutions  $u$  of (2.1) for a given parameter  $\gamma = (\xi, \phi_1, \phi_2, r, u_0) \in \Gamma$ . Defining  $u(t) \equiv u(t, \cdot) \in L_2(0,1)$ , we may rewrite (2.1) as an initial value problem in  $u$ ,

$$(2.4) \quad \begin{cases} u_t = A(q)u(t) + F(t;r), & t \in (0,T), \\ u(0) = u_0. \end{cases}$$

Here  $q = \phi_1 + H_\xi \phi_2$ ,  $F(t;r) = f(t, \cdot; r(\cdot))$  and the operator  $A(q)$  is defined by  $A(q)\psi = \mathcal{D}(q\mathcal{D}\psi)$  for  $\psi \in \text{dom}A(q) = V_q$ , where

$$V_q = \{\psi \in H_0^1(0,1) \mid q\mathcal{D}\psi \in H^1(0,1)\}$$

(throughout we shall use  $\mathcal{D}$  to denote the spatial differentiation operator  $\frac{\partial}{\partial x}$ ).

We note that it would be natural, given the discontinuities in coefficients  $\rho$  and  $q$ , to consider a weak form of (2.1) in order to relax restrictions on both solutions and parameters. For this particular problem, however, we shall insist that solutions  $u$  satisfy a continuity equation

$$(2.5) \quad (q\mathcal{D}u)(\xi^-) = (q\mathcal{D}u)(\xi^+),$$

which represents continuity of stress across a transition point,  $\xi$ , between

distinct spatial regions (layers of porous media, for example). Given this condition on  $u$  (which implies that  $qDu$  must be sufficiently regular to ensure that point evaluations are meaningful), there is reason for seeking at least strong solutions to (2.1). We have one additional comment along these lines from a numerical point of view: It is of interest to note that we experienced no difficulties in applying our state variable approximation ideas to the forward problem (i.e., the problem of integrating (2.1), or a weak form of (2.1), for a known value of  $\gamma$ ) for particular problems where the true state  $u$  did not satisfy (2.5). It was not until we turned to the numerical solution of the inverse problem (and, in particular, the problem of estimating  $\xi$  itself) that it became evident that one could not expect to estimate  $\xi$  unless the (physically meaningful) continuity equation (2.5) was satisfied by solutions of (2.1). Therefore, there seems to be little justification in considering weak solutions of (2.1) in the context of the parameter estimation problem.

Our first result is a statement of existence and uniqueness of solutions of (2.1); in addition, we indicate certain regularity properties of solutions that will be useful in later calculations.

Theorem 2.1. Let  $\gamma = (\xi, \phi_1, \phi_2, r, u_0)$  be given in  $\Gamma$  and let  $q = \phi_1 + H_\xi \phi_2$ . There exists a unique (classical) solution  $u$  to (2.1) with the property that  $u(t) \in V_q$  for any  $t > 0$ . In addition, if  $u_0 \in V_q$ , then the map  $t \rightarrow A(q)u(t)$  is in  $C([0, T]; L_2(0, 1))$ .

Proof: It suffices to show that  $A(q)$  generates an analytic semigroup on  $L_2(0, 1)$ . From this we may guarantee existence of a unique solution  $u$  to (2.1) (see Corollary 3.3, p. 113 of [27]); we may then apply the well-known smoothing properties of analytic semigroups (see, for example, Theorem 3.5, p. 114, of [27]),

along with hypothesis (H2), to obtain the statement of the remainder of the theorem.

We first show that  $A(q)$  is densely defined and self-adjoint. To this end we note that  $\tilde{V}_q \equiv \{\psi \in L_2(0,1) \mid \psi \in H_0^2(0,\xi), \psi \in H_0^2(\xi,1)\}$  satisfies  $\tilde{V}_q \subseteq V_q = \text{dom}A(q)$  (since  $qD\psi(\xi^-) = qD\psi(\xi^+) = 0$ ); it is easy to argue that  $\tilde{V}_q$  is dense in  $L_2(0,1)$  if one uses the fact that  $H_0^2(0,\xi)$  and  $H_0^2(\xi,1)$  are dense in  $L_2(0,\xi)$  and  $L_2(\xi,1)$ , respectively. Using an integration by parts, it is not difficult to show that  $A(q)$  is symmetric. To demonstrate that  $A(q)$  is self-adjoint, it suffices to show [28; Theorem 13.11] that  $\text{Range } A(q) = L_2(0,1)$ ; that is, for  $g \in L_2(0,1)$ , it is sufficient to verify the existence of a solution to

$$(2.6) \quad A(q)\psi = g$$

that satisfies  $\psi \in V_q$ . Because one may readily see this is true (using standard theory for two-point ordinary differential equation boundary value problems--see, for example, Theorem 8.3 of [19]), it therefore follows that  $A(q)$  is self-adjoint. In addition,  $A(q)$  is dissipative (since  $\langle A(q)\psi, \psi \rangle \leq -m|D\psi|^2 \leq 0$  for all  $\psi \in V_q$ ) so that a corollary of the Lumer Phillips theorem [27; p. 15] may be invoked to argue that  $A(q)$  generates a  $C_0$  semigroup of contractions on  $L_2(0,1)$ . Finally, since  $\sigma(A(q)) \subseteq (-\infty, 0)$ , we may apply standard semigroup theory [20; Theorem 7.12, p. 82], [27; p. 61] to conclude that  $A(q)$  generates an analytic semigroup (analytic on the sector  $\{\lambda \in \mathbb{C} \mid \lambda \neq 0, |\arg \lambda| < \frac{\pi}{2}\}$ ). The proof of the theorem is complete.

It is useful to note that if  $u$  is a solution of (2.1) then  $u$  also satisfies (2.1) in a weak sense; i.e.,  $u$  satisfies

$$(2.7) \quad \begin{cases} \langle u_t(t), v \rangle = - \langle q \mathcal{D}u(t), \mathcal{D}v \rangle + \langle F(t; r), v \rangle, & t \in (0, T) \\ u(0) = u_0 \end{cases}$$

for any  $v \in H_0^1(0,1)$ . Using this formulation we may argue the continuous dependence of solutions on (possibly unknown) initial data. In fact, we may actually show that the map  $\gamma \rightarrow u(t; \gamma) : \Gamma \rightarrow L_2(0,1)$  is continuous, uniform in  $t$ , so that we are guaranteed the existence of a minimizer for  $J$  over the (compact) set  $\Gamma$ . As we also establish existence of an "optimal" parameter in Theorem 3.3 (and because we do not need the continuity of  $\gamma \rightarrow u(t; \gamma)$  to make any of the arguments given below), we state and prove only the result that follows.

Corollary 2.1. The mapping  $u_0 \rightarrow u(t; \xi, \phi_1, \phi_2, r, u_0) : L_2(0,1) \rightarrow L_2(0,1)$  is continuous, uniform in  $(\xi, \phi_1, \phi_2, r) \in \mathcal{S}$  and  $t \in (0, T)$ .

Proof: Let  $u_0, u'_0$  be given in  $L_2(0,1)$  and define  $u(t) \equiv u(t; \xi, \phi_1, \phi_2, r, u_0)$ ,  $u'(t) \equiv u(t; \xi, \phi_1, \phi_2, r, u'_0)$ . Using (2.7), we find that, for  $t > 0$ ,

$$\langle u_t(t) - u'_t(t), v \rangle = - \langle q \mathcal{D}(u(t) - u'(t)), \mathcal{D}v \rangle$$

for any choice of test function  $v$  in  $H_0^1$ ; setting  $v = u(t) - u'(t) \in H_0^1$  (from Theorem 2.1) we obtain

$$\frac{1}{2} \frac{d}{dt} |u(t) - u'(t)|^2 + m |\mathcal{D}(u(t) - u'(t))|^2 \leq 0.$$

In fact, using the Rayleigh-Ritz inequality [29; p. 5], it follows that

$$\frac{1}{2} \frac{d}{dt} |u(t) - u'(t)|^2 + m\pi^2 |u(t) - u'(t)|^2 \leq 0$$

so that an application of the Gronwall inequality yields



$$|u(t) - u'(t)|^2 \leq e^{-2m\pi^2 t} |u_0 - u'_0|^2 .$$

Continuous dependence of  $u(t;s,u_0)$  on  $u_0$ , uniform in  $s \in \mathcal{S}$ , thus obtains.

We are ready to consider the problem of unknown parameters, in the context of the parameter estimation problems defined in this section. We note that the problem of estimating an optimal parameter  $\gamma^* \in \Gamma$  must be combined with schemes for solving (2.1); i.e., we must consider state variable approximation as well as the problem of finding finite-dimensional approximations for unknown functional parameters. In the sections that follow, we develop both state and parameter approximation techniques with the goal of solving these problems.

### 3. A Spline-Based Approximation Scheme

Standard numerical optimization schemes applied to the problem of minimizing  $J$  (or  $\tilde{J}$ ) over  $r$  typically generate a minimizing sequence of parameter iterates, starting from an initial guess,  $\gamma^0$ . However, schemes of this type generally require that  $u(\gamma)$  (the solution of (2.1)) be evaluated as the parameter  $\gamma$  is updated; it is therefore desirable to combine estimation of an optimal parameter  $\gamma^*$  with approximation techniques for solving (2.1). With this goal in mind, we describe a spline-based state/parameter approximation scheme in the same spirit of the ideas found in [5], [9], [10], [12], [14], [23] to name a few of the related references in this area for (continuous coefficient) parabolic problems.

The convergence arguments developed below are similar to standard variational-type estimates often used in association with finite element approximations (see, for example, [29], p. 129) although the estimates given here are complicated somewhat by the presence of unknown parameters. This variational approach was taken in [10] and [11] for the problem of estimating continuous coefficients in parabolic systems; we require a somewhat different treatment here primarily due to the fact that we allow discontinuous coefficients, where the points of discontinuity are unknown (necessitating parameter-dependent approximation spaces  $X^N(q)$ ). Thus, an interesting aspect of our approach (and often a source of difficulties) involves the fact that our approximation spaces change with every choice of parameter iterate. We note that although the theoretical problems are quite different, our construction of approximating spaces  $X^N(q)$  is somewhat similar to the ideas found in [2], [13], [16]; there the problem was to estimate unknown delays appearing in functional differential equations (there is a correspondence between our treatment of an unknown point of discontinuity and the approach taken in those references to handle an unknown

delay, at least from the standpoint of numerical approximation schemes). We turn now to a precise statement of the approximation scheme under consideration.

For any  $\gamma = (\xi, \phi_1, \phi_2, r, u_0) \in \Gamma$ , we construct parameter-dependent spaces and operators as follows: For  $q = \phi_1 + H_\xi \phi_2$  and  $N = 1, 2, \dots$ , we define  $X^N(q) = \text{span}\{B_i^N(q), i = 1, \dots, 2N-1\}$ , where  $B_i^N(q)$  denotes the  $i^{\text{th}}$  continuous piecewise-linear B-spline basis element (satisfying homogeneous boundary conditions) with knots at  $\{x_k^N(q), k = 0, \dots, 2N\}$ . Here  $x_k^N(q) = k\xi/N$ ,  $k = 0, \dots, N$ , and  $x_k^N(q) = \xi + (k-N)(1-\xi)/N$ , for  $k = N+1, \dots, 2N$ . The piecewise linear elements are characterized by  $B_i^N(q)(x_k^N) = \delta_{ik}$  for  $i, k = 1, \dots, 2N-1$  ( $B_i^N(q)(0) = B_i^N(q)(1) = 0$ ); see Figures 1 - 3.

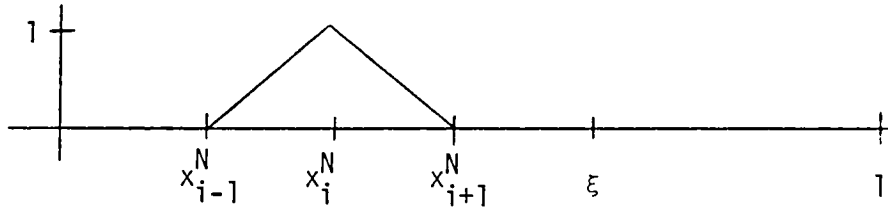


Figure 1.  $B_i^N$ ,  $i = 1, \dots, N-1$ .

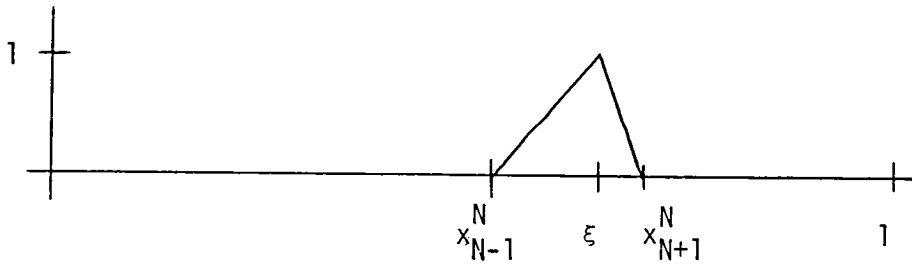


Figure 2.  $B_N^N$ .



Figure 3.  $B_i^N$ ,  $i = N+1, \dots, 2N-1$ .

We remark that in general, for  $\gamma, \tilde{\gamma} \in \Gamma$  and  $q = \phi_1 + H_\xi \phi_2$ ,  $\tilde{q} = \tilde{\phi}_1 + H_{\tilde{\xi}} \tilde{\phi}_2$ , we do not have  $X^N(q) \subseteq X^N(\tilde{q})$  or  $X^N(\tilde{q}) \subseteq X^N(q)$ , nor do we have  $X^N(q) \subseteq V_q$  (note that although an element  $\psi^N \in X^N(q)$  does have a discontinuity in its first derivative at  $\xi$ ,  $\psi^N$  does not satisfy the continuity equation (2.5)

associated with  $q$ . The approximation spaces  $X^N(q)$  are chosen so that the resulting parameter estimation algorithm enjoys a number of computational advantages, especially when  $\xi$  is unknown. We take a general Galerkin approach to define approximating state systems and then obtain convergence findings by working directly with the weak form of these equations. As an alternative approach to that taken here one could define approximating operators  $A^N(q)$  for  $A(q)$  and investigate the sense in which  $A^N(q)$  "converges" to  $A(q)$  (see, for example, Example 2.2 of [22], or [23], for approximating operators that might be used in this context).

For  $\gamma \in \Gamma$  fixed and  $N = 1, 2, \dots$ , we seek an approximation to  $u(t; \gamma)$  of the form  $u^N(t; \gamma) = \sum_{i=1}^{2N-1} w_i^N(t; \gamma) B_i^N(q)$ , where the "Fourier coefficients"  $w_i^N$  are determined via the system of ordinary differential equations (ODE),

$$(3.1) \quad \begin{cases} \langle u_t^N(t; \gamma), B_i^N(q) \rangle = - \langle q \mathcal{D} u^N(t; \gamma), \mathcal{D} B_i^N(q) \rangle + \langle F(t; r), B_i^N(q) \rangle, & t \in (0, T), \\ \langle u^N(0; \gamma), B_i^N(q) \rangle = \langle u_0, B_i^N(q) \rangle, \end{cases}$$

for  $i = 1, \dots, 2N-1$ . Alternatively,  $u^N$  satisfies

$$(3.2) \quad \begin{cases} \langle u_t^N(t; \gamma), v \rangle = - \langle q \mathcal{D} u^N(t; \gamma), \mathcal{D} v \rangle + \langle F(t; r), v \rangle, & t \in (0, T), \\ u^N(0; \gamma) = P^N(q) u_0 \end{cases}$$

for all  $v \in X^N(q)$ ; here  $P^N(q) : L_2(0, 1) \rightarrow X^N(q)$  denotes the orthogonal projection (with respect to the usual  $L_2$  topology) along  $X^N(q)^\perp$ . Associated with (3.1) is an approximate estimation problem, namely that of finding  $\bar{\gamma}^N \in \Gamma$  that minimizes

$$(3.3) \quad J^N(\gamma) = \sum_{i=1}^n \int_0^1 |\mathcal{C}(t_i, x; \gamma) u^N(t_i; \gamma)(x) - \hat{u}_i(x)|^2 dx$$

over  $\Gamma$ , where  $u^N(\gamma)$  is the solution of (3.1) corresponding to  $\gamma \in \Gamma$ .

Our initial findings concerning the  $N^{\text{th}}$  approximate problem (3.1), (3.3) are immediate consequences of the fact that (3.1) is an ODE on  $X^N(q)$  and that the basis elements  $B_i^N(q)$  (and their spatial derivatives) are continuous in  $\xi$  (see (4.1) for a more explicit matrix representation of (3.1)). We shall defer a more detailed examination of (3.1) (and the implementation of the estimation scheme associated with (3.1), (3.3)), until Section 4 where our numerical findings are summarized.

**Theorem 3.1.** For each  $N$  and any  $\gamma \in \Gamma$ , there exists a unique solution  $u^N(\gamma)$  of (3.1),  $u^N(t; \gamma) \in X^N(q)$ . In addition, the mapping  $\gamma \mapsto u^N(t; \gamma) : \Gamma \rightarrow L_2(0, 1)$  is continuous for each  $t \in (0, T)$ .

Corollary 3.1. For each  $N$ , there exists a solution  $\bar{\gamma}^N \in \Gamma$  for the problem of minimizing  $J^N$  over  $\Gamma$ .

Finally, a simple modification of the proof of Corollary 2.1 yields a similar statement concerning the continuous dependence of  $u^N(t)$  on  $u_0$ :

Corollary 3.2. The mapping  $u_0 \rightarrow u^N(t; (s, u_0)) : L_2(0,1) \rightarrow L_2(0,1)$  is continuous, uniform in  $N$ ,  $s = (\xi, \phi_1, \phi_2, r) \in \mathcal{S}$ , and  $t \in (0,1)$ .

An essential step in the process of correlating state variable approximation with the problem of estimating an optimal parameter  $\gamma^* \in \Gamma$  (for the original parameter identification problem) is the establishment of the convergence of  $u^N(t; \gamma^N)$  to  $u(t; \gamma)$  for any sequence  $\{\gamma^N\}$  in  $\Gamma$  that converges to  $\gamma \in \Gamma$ . We shall clarify the need for arguments of this type in the proof of Theorem 3.3. To facilitate steps in this direction, we shall first establish linear spline estimates, the proof of which are in the spirit of [29; pp. 16-17, 78]. In what follows we assume that  $\{\gamma^N\}$  is given in  $\Gamma$ ,  $\gamma^N = (\xi^N, \phi_1^N, \phi_2^N, r^N, u_0^N)$ , with  $\gamma^N \rightarrow \bar{\gamma} = (\bar{\xi}, \bar{\phi}_1, \bar{\phi}_2, \bar{r}, \bar{u}_0) \in \Gamma$  (in the usual product topology on  $\Gamma$ ); in addition, we assume there exists  $\delta$  such that  $0 < \delta \leq \bar{\xi} \leq 1 - \delta < 1$  (and, in the case of multiple discontinuities,  $|\bar{\xi}_k - \bar{\xi}_{k-1}| \geq \delta > 0$ ,  $k = 1, \dots, \mu$ ). Given  $q^N = \phi_1^N + H_{\xi^N} \phi_2^N$ , we shall henceforth simplify notation and abbreviate  $p^N \equiv p^N(q^N)$ ,  $x^N \equiv x^N(q^N)$ , and  $x_k^N \equiv x_k^N(q^N)$ ,  $k = 0, \dots, 2N$ .

Lemma 3.1. Let  $\psi$  be given in  $V_{\bar{q}}$ , where  $\bar{q} = \bar{\phi}_1 + H_{\bar{\xi}} \bar{\phi}_2$ . There exist constants  $c_1$  and  $c_2$ , independent of  $N$ , such that

$$(3.4) \quad |\psi - p^N \psi| \leq c_1 N^{-2} |A(\bar{q})\psi|$$

and, for  $N$  sufficiently large,

$$(3.5) \quad |\mathfrak{D}(\psi - \mathcal{P}^N \psi)| \leq c_2 N^{-1} |A(\bar{q})\psi|.$$

Proof. We shall denote by  $\mathcal{I}^N \psi$  the linear interpolant of  $\psi$ , with knots at  $x_k^N$ ,  $k = 0, \dots, 2N$ ; that is,  $\mathcal{I}^N \psi(x_k^N) = \psi(x_k^N)$ ,  $k = 0, \dots, 2N$ . We find that

$$|\mathfrak{D}(\psi - \mathcal{I}^N \psi)|^2 = \langle \mathfrak{D}\psi, \mathfrak{D}(\psi - \mathcal{I}^N \psi) \rangle - \langle \mathfrak{D}(\mathcal{I}^N \psi), \mathfrak{D}(\psi - \mathcal{I}^N \psi) \rangle$$

where

$$\begin{aligned} \langle \mathfrak{D}(\mathcal{I}^N \psi), \mathfrak{D}(\psi - \mathcal{I}^N \psi) \rangle &= \sum_{k=1}^{2N} \int_{x_{k-1}^N}^{x_k^N} \mathfrak{D}(\mathcal{I}^N \psi) \mathfrak{D}(\psi - \mathcal{I}^N \psi) \\ &= - \sum_{k=1}^{2N} \int_{x_{k-1}^N}^{x_k^N} \mathfrak{D}^2(\mathcal{I}^N \psi) (\psi - \mathcal{I}^N \psi) \\ &= 0. \end{aligned}$$

Thus,

$$\begin{aligned} |\mathfrak{D}(\psi - \mathcal{I}^N \psi)|^2 &= \langle \mathfrak{D}\psi, \mathfrak{D}(\psi - \mathcal{I}^N \psi) \rangle \\ &\leq \frac{1}{m} \sum_{k=1}^{2N} \int_{x_{k-1}^N}^{x_k^N} \bar{q} \mathfrak{D}\psi \mathfrak{D}(\psi - \mathcal{I}^N \psi) \\ &= - \frac{1}{m} \sum_{k=1}^{2N} \int_{x_{k-1}^N}^{x_k^N} \mathfrak{D}(\bar{q} \mathfrak{D}\psi) (\psi - \mathcal{I}^N \psi) \\ &\leq \frac{1}{m} |\mathfrak{D}(\bar{q} \mathfrak{D}\psi)| |\mathcal{I}^N \psi - \psi| \\ &\leq (m\pi N)^{-1} |A(\bar{q})\psi| |\mathfrak{D}(\psi - \mathcal{I}^N \psi)| \end{aligned}$$

where we have used (2.16) of [29; p. 17] in the last estimate. It therefore follows that

$$(3.6) \quad |\mathfrak{D}(\psi - \mathcal{I}^N \psi)| \leq (m\pi N)^{-1} |\mathfrak{A}(\bar{q})\psi|$$

and, again using (2.16) from [29],

$$(3.7) \quad |\mathcal{I}^N_{\psi-\psi}| \leq (m\pi^2 N^2)^{-1} |\mathfrak{A}(\bar{q})\psi|.$$

To establish (3.4), we use properties of the projection  $P^N$  to note that

$$|\psi - P^N \psi| \leq |\psi - \mathcal{I}^N \psi| \leq c_1 N^{-2} |\mathfrak{A}(\bar{q})\psi|.$$

Finally,  $|\mathfrak{D}(\psi - P^N \psi)| \leq |\mathfrak{D}(\psi - \mathcal{I}^N \psi)| + |\mathfrak{D}(\mathcal{I}^N_{\psi-P^N \psi})|$ , where an application of the Schmidt inequality [29; p. 7] yields (for  $N$  sufficiently large)

$$\begin{aligned} |\mathfrak{D}(\mathcal{I}^N_{\psi-P^N \psi})|^2 &= \sum_{k=1}^{2N} \int_{x_{k-1}^N}^{x_k^N} |\mathfrak{D}(\mathcal{I}^N_{\psi-P^N \psi})|^2 \\ &\leq 12 \sum_{k=1}^{2N} (x_k^N - x_{k-1}^N)^{-2} \int_{x_{k-1}^N}^{x_k^N} |\mathcal{I}^N_{\psi-P^N \psi}|^2 \\ &\leq 12 (N/\delta)^2 |\mathcal{I}^N_{\psi-P^N \psi}|^2 \end{aligned}$$

so that

$$\begin{aligned} |\mathfrak{D}(\psi - P^N \psi)| &\leq |\mathfrak{D}(\psi - \mathcal{I}^N \psi)| + 2\sqrt{3}(N/\delta) |\mathcal{I}^N_{\psi-P^N \psi}| \\ &\leq |\mathfrak{D}(\psi - \mathcal{I}^N \psi)| + 2\sqrt{3}(N/\delta) \{ |\mathcal{I}^N_{\psi-\psi}| + |\psi - P^N \psi| \} \\ &\leq (m\pi N)^{-1} |\mathfrak{A}(\bar{q})\psi| + 4\sqrt{3}(N/\delta) c_1 N^{-2} |\mathfrak{A}(\bar{q})\psi| \\ &\leq c_2 N^{-1} |\mathfrak{A}(\bar{q})\psi| \end{aligned}$$

for an appropriate choice of the constant  $c_2$ .



We may now use the linear spline estimates derived in Lemma 3.1 to establish a preliminary convergence result.

Lemma 3.2. Suppose  $\{\gamma^N\}$  is given in  $\Gamma$  such that  $\gamma^N \rightarrow \bar{\gamma} = (\bar{\xi}, \bar{\phi}_1, \bar{\phi}_2, \bar{r}, \bar{u}_0) \in \Gamma$ . Assume, in addition, that  $\bar{u}_0 \in V_{\bar{q}}$ , where  $\bar{q} = \bar{\phi}_1 + H_{\bar{\xi}} \bar{\phi}_2$ . Then, for every  $t \in (0, T)$ ,

$$u^N(t; \gamma^N) \rightarrow u(t; \bar{\gamma}) \text{ in } L_2(0, 1)$$

as  $N \rightarrow \infty$  (where  $u^N$  is the solution of (3.1) associated with  $\gamma^N$  and  $u$  is the solution of (2.1) associated with  $\bar{\gamma}$ ).

Proof. Let  $u(t) \equiv u(t; \bar{\gamma})$ ,  $u^N(t) \equiv u^N(t; \gamma^N)$ . Then

$$|u^N(t) - u(t)| \leq |u^N(t) - P^N u(t)| + |P^N u(t) - u(t)|$$

where the second term is  $O(N^{-2})$  from (3.4) ( $u(t) \in V_{\bar{q}}$  for  $t \in (0, T)$ ; Theorem 2.1). To consider the first term, we note that solutions  $u(t)$ ,  $u^N(t)$  of (2.1), (3.1), respectively, satisfy (2.7) and (3.2) for any  $v \in X^N \subseteq H_0^1(0, 1)$ . Using these equations it is easy to see that

$$\begin{aligned} & \langle u_t^N(t), v \rangle + \langle q^N \mathcal{D} u^N(t), \mathcal{D} v \rangle - \langle F^N(t), v \rangle - \left( \langle \frac{d}{dt} P^N u(t), v \rangle + \langle q^N \mathcal{D} P^N u(t), \mathcal{D} v \rangle \right) \\ &= 0 - \left( \langle \frac{d}{dt} P^N u(t), v \rangle + \langle q^N \mathcal{D} P^N u(t), \mathcal{D} v \rangle \right) \\ &+ \langle u_t(t), v \rangle + \langle \bar{q} \mathcal{D} u(t), \mathcal{D} v \rangle - \langle F(t), v \rangle \end{aligned}$$

where  $F(t) \equiv F(t; \bar{r})$ ,  $F^N(t) \equiv F(t; r^N)$ , and  $q^N = \phi_1^N + H_{\xi^N} \phi_2^N$ . It thus follows that

$$\begin{aligned} & \langle \frac{d}{dt} (u^N(t) - P^N u(t)), v \rangle = - \langle q^N \mathcal{D} (u^N(t) - P^N u(t)), \mathcal{D} v \rangle \\ &+ \langle \frac{d}{dt} (u(t) - P^N u(t)), v \rangle + \langle \bar{q} \mathcal{D} u(t) - q^N \mathcal{D} P^N u(t), \mathcal{D} v \rangle + \langle F^N(t) - F(t), v \rangle. \end{aligned}$$

Letting  $v = u^N(t) - p^N_u(t) \in X^N$ , we argue that

$$\begin{aligned}
\frac{1}{2} \frac{d}{dt} |u^N(t) - p^N_u(t)|^2 &\leq -m |\mathfrak{D}(u^N(t) - p^N_u(t))|^2 \\
&\quad + \left| \frac{d}{dt} (u(t) - p^N_u(t)) \right| |u^N(t) - p^N_u(t)| \\
&\quad + \left( (2m)^{-\frac{1}{2}} |\bar{q}\mathfrak{D}u(t) - q^N\mathfrak{D}(p^N_u(t))| \right) \left( (2m)^{\frac{1}{2}} |\mathfrak{D}(u^N(t) - p^N_u(t))| \right) \\
&\quad + |F^N(t) - F(t)| |u^N(t) - p^N_u(t)| \\
&\leq \frac{1}{2} \left| \frac{d}{dt} (u(t) - p^N_u(t)) \right|^2 + \frac{1}{4m} |\bar{q}\mathfrak{D}u(t) - q^N\mathfrak{D}(p^N_u(t))|^2 \\
&\quad + \frac{1}{2} |F^N(t) - F(t)|^2 + |u^N(t) - p^N_u(t)|^2,
\end{aligned}$$

where we have repeatedly used the inequality  $ab \leq \frac{1}{2}(a^2 + b^2)$ . Defining  $w^N(t) \equiv |u^N(t) - p^N_u(t)|^2$ , the above estimates reduce to

$$\begin{aligned}
\dot{w}^N(t) - 2w^N(t) &\leq \left| \frac{d}{dt} (u(t) - p^N_u(t)) \right|^2 + \frac{1}{2m} |\bar{q}\mathfrak{D}u(t) - q^N\mathfrak{D}(p^N_u(t))|^2 \\
&\quad + |F^N(t) - F(t)|^2,
\end{aligned}$$

so that an application of the Gronwall inequality yields

$$w^N(t) = |u^N(t) - p^N_u(t)|^2 \leq e^{2T} \{\tau_1^N + \tau_2^N + \tau_3^N + \tau_4^N\}$$

where

$$\begin{aligned}
\tau_1^N &= |u^N(0, \gamma^N) - p^N_u(0; \bar{\gamma})|^2 \\
\tau_2^N &= \int_0^T \left| \frac{d}{dt} (u(s; \bar{\gamma}) - p^N_u(s; \bar{\gamma})) \right|^2 ds \\
\tau_3^N &= \frac{1}{2m} \int_0^T |\bar{q}\mathfrak{D}u(s, \bar{\gamma}) - q^N\mathfrak{D}(p^N_u(s; \bar{\gamma}))|^2 ds \\
\tau_4^N &= \int_0^T |F(s; r^N) - F(s; \bar{r})|^2 ds.
\end{aligned}$$

From hypothesis (H3), it follows that  $\tau_4^N \rightarrow 0$  as  $N \rightarrow \infty$ . We are also able to argue the convergence of  $\tau_1^N$  to 0 since  $\tau_1^N = |p^N u_0^N - p^N \bar{u}_0^N| \leq |u_0^N - \bar{u}_0^N|$ .

To consider  $\tau_2^N$ , it is useful to note that, for  $v$  fixed in  $X^N$ , the function defined by  $g^N(t) \equiv \langle u(t) - p^N u(t), v \rangle$  is identically zero (from the definition of  $p^N$ ) so that  $0 = g_t^N(t) = \langle u_t(t) - \frac{d}{dt} p^N u(t), v \rangle$ . But this is true for every  $v \in X^N$  and all  $t \in (0, T)$ , so it must follow that  $\frac{d}{dt} p^N u(t) = p^N u_t(t)$ . Thus,

$$\tau_2^N = \int_0^T |u_t(s) - p^N u_t(s)|^2 ds$$

where, for each  $s \in (0, T)$ , the integrand converges to zero as  $N \rightarrow \infty$ ; this claim may be verified using (3.4) and the fact that  $u_t(s) \in L_2(0, 1)$  and  $V_q$  is dense in  $L_2(0, 1)$ . The integrand is dominated by  $2|u_t(s)|^2$  where  $s \rightarrow u_t(s)$  is in  $L_2((0, T), L_2(0, 1))$  (Theorem 2.1). It follows that  $\tau_2^N \rightarrow 0$ .

Finally, for  $N$  sufficiently large,

$$\begin{aligned} m\tau_3^N &\leq \int_0^T |(\bar{q} - q^N)Du(s)|^2 + \int_0^T |q^N D(u(s) - p^N u(s))|^2 \\ &\leq 2 \int_0^T |(\bar{\phi}_1 - \phi_1^N)Du(s)|^2 + 4 \int_0^T |H_{\bar{\xi}}(\bar{\phi}_2 - \phi_2^N)Du(s)|^2 + 4 \int_0^T |(H_{\bar{\xi}} - H_{\xi^N})\phi_2^N Du(s)|^2 \\ &\quad + \bar{m}^2 \int_0^T |D(u(s) - p^N u(s))|^2 \\ &\leq 4 (|\bar{\phi}_1 - \phi_1^N|_{\infty}^2 + |\bar{\phi}_2 - \phi_2^N|_{\infty}^2) \int_0^T |Du(s)|^2 + 2\bar{m}^2 |\bar{\xi} - \xi^N| \int_0^T |Du(s)|^2 \\ &\quad + (\bar{m}c_2 N^{-1})^2 \int_0^T |A(\bar{q})u(s)|^2 \end{aligned}$$

where we have used (3.5) in the last inequality. Further,  $\int_0^T |Du(s)|^2 = \int_0^T \langle Du(s), Du(s) \rangle \leq m^{-1} \int_0^T \langle \bar{q} Du(s), Du(s) \rangle \leq m^{-1} \int_0^T |A(\bar{q})u(s)| |u(s)| < \infty$

(since  $\bar{u}_0 \in V_{\bar{q}}$ ) so that  $\tau_3^N \rightarrow 0$  as  $N \rightarrow \infty$ .

Finally, it is possible to lift the requirement that  $\bar{u}_0$  belongs to  $V_{\bar{q}}$  and to prove a more general convergence theorem.

Theorem 3.2. Assume  $\{\gamma^N\}$  is given in  $\Gamma$ ,  $\gamma^N = (\xi^N, \phi_1^N, \phi_2^N, r^N, u_0^N)$ , such that  $\gamma^N \rightarrow \bar{\gamma} = (\bar{\xi}, \bar{\phi}_1, \bar{\phi}_2, \bar{r}, \bar{u}_0) \in \Gamma$ . Then, for each  $t \in (0, T)$ ,

$$u^N(t; \gamma^N) \rightarrow u(t; \bar{\gamma}) \text{ in } L_2(0, 1)$$

as  $N \rightarrow \infty$ .

Proof. Let  $\epsilon > 0$  be given and define  $\bar{q} = \bar{\phi}_1 + H_{\bar{\xi}} \bar{\phi}_2$ ,  $s^N = (\xi^N, \phi_1^N, \phi_2^N, r^N)$ , and  $\bar{s} = (\bar{\xi}, \bar{\phi}_1, \bar{\phi}_2, \bar{r})$ . Since  $V_{\bar{q}}$  is dense in  $L_2(0, 1)$ , there exists  $\psi \in V_{\bar{q}}$  sufficiently close to  $\bar{u}_0 \in L_2(0, 1)$  so that we may argue that

$$\begin{aligned} |u^N(t; \gamma^N) - u(t; \bar{\gamma})| &\leq |u^N(t; (s^N, u_0^N)) - u^N(t; (s^N, \psi))| \\ &\quad + |u^N(t; (s^N, \psi)) - u(t; (\bar{s}, \psi))| \\ &\quad + |u(t; (\bar{s}, \psi)) - u(t; (\bar{s}, \bar{u}_0))| \\ &< \epsilon, \end{aligned}$$

for  $N$  sufficiently large.

Here we have used the continuous dependence of  $u^N$ ,  $u$  on initial conditions (uniform in  $s^N$ ,  $\bar{s}$ , and  $N$ ), the inequality  $|u_0^N - \psi| \leq |u_0^N - \bar{u}_0| + |\bar{u}_0 - \psi|$ , and the findings in Lemma 3.2. The proof of the theorem thus obtains.

To this point, we have focused on state variable convergence (of  $u^N$  to  $u$ ) once the convergence of any sequence of parameters has been guaranteed. In reality, however, we have yet to establish whether any sequence of solutions  $\{\bar{\gamma}^N\}$  for the approximating estimation problems is indeed convergent; even then we have no assurance that the limiting parameter  $\bar{\gamma}$  is in fact a solution to the original parameter identification problem. In our next theorem we consider this problem and indicate when an approximate estimation problem may be used numerically to compute an approximate solution for the original problem. The proof of the theorem is similar to ideas found in [12], [14].

Theorem 3.3. For each  $N$  let  $\bar{\gamma}^N$  denote a solution for the problem of minimizing  $J^N$  over  $\Gamma$ . There exists  $\gamma^* \in \Gamma$  and a subsequence  $\{\bar{\gamma}^{N_k}\}$  of  $\{\bar{\gamma}^N\}$  such that

- (i)  $\bar{\gamma}^{N_k} \rightarrow \gamma^*$  in the product topology on  $\Gamma$ ,
- (ii)  $u^{N_k}(t; \bar{\gamma}^{N_k}) \rightarrow u(t; \gamma^*)$  for each  $t \in (0, T)$ ,
- (iii)  $J^{N_k}(\bar{\gamma}^{N_k}) \rightarrow J(\gamma^*)$ , and,
- (iv)  $\gamma^*$  is a solution to the original parameter estimation problem, namely that of minimizing  $J(\gamma)$  over  $\Gamma$ .

Proof. Parts (i) - (iii) are immediate consequences of hypothesis (H1), and Theorem 3.2. To prove part (iv) it suffices to note that

$$\begin{aligned}
 J(\gamma^*) &= \lim_{N_k \rightarrow \infty} J^{N_k}(\bar{\gamma}^{N_k}) \\
 &\leq \lim_{N_k \rightarrow \infty} J^{N_k}(\gamma) \\
 &= J(\gamma)
 \end{aligned}$$

for any  $\gamma \in \Gamma$  ( $\bar{\gamma}^{N_k}$  is a minimizer for  $J^{N_k}$  over  $\Gamma$ ), so that  $\gamma^*$  is a solution for the problem of minimizing  $J$  over  $\Gamma$ .

**Remark 3.1.** Although it has been assumed in our discussion so far that  $\rho \equiv 1$ , there are no difficulties associated with extending our ideas to more general  $\rho$  ( $\rho = \kappa_1 + H_\xi \kappa_2$ ,  $(\kappa_1, \kappa_2) \in K \equiv \{(\kappa_1, \kappa_2) \in C[0,1] \times C[0,1] \mid 0 < m \leq \kappa_i \leq \bar{m}\}$ ). For example, it is easy to see that the arguments used in the proof of Theorem 2.1 change very little if, instead of  $A(q)$  and the  $L_2$  inner product, one considers the operator  $B(\rho, q) \equiv \frac{1}{\rho} A(q)$  and the weighted  $L_2$  inner product defined by  $\langle \psi, \zeta \rangle_\rho = \int_0^1 \rho \psi \zeta$ . For each  $N$  and  $\gamma^N = (\xi^N, \phi_1^N, \phi_2^N, r^N, u_0^N, \kappa_1^N, \kappa_2^N)$  given in  $\Gamma$  ( $\Gamma$  a compact subset of  $\mathcal{S} \times L_2(0,1) \times K$ ), we may define approximating equations in the variable  $u^N(t) \in X^N(q^N)$  by

$$\begin{cases} \langle \rho^N u_t^N(t), v \rangle = - \langle q^N D u^N(t), D v \rangle + \langle F(t; r^N), v \rangle & , v \in X^N(q^N) \\ u^N(0) = p^N u_0^N, \end{cases}$$

where  $\rho^N = \kappa_1^N + H_{\xi^N} \kappa_2^N$  and  $q^N = \phi_1^N + H_{\xi^N} \phi_2^N$ . To establish convergence of  $u^N$  to  $u$  (as  $\gamma^N \rightarrow \bar{\gamma}$ ) one may make a simple modification in the proof of Lemma 3.2 to argue that

$$\begin{aligned} m |u^N(t) - p^N u(t)|^2 &\leq |(\rho^N)^{1/2} (u^N(t) - p^N u(t))|^2 \\ &\leq ce^{2T\{\tau_1^N + \tau_3^N + \tau_4^N + |\bar{\rho} u_t - \rho^N p^N u_t|^2\}} \end{aligned}$$

where the  $\tau_i^N$  are unchanged from that lemma and the last term converges to zero. Because all other estimates remain unchanged, we are able to derive an analog to Theorem 3.3, i.e., we are able to treat the case of a general (possibly unknown) parameter  $\rho$ .

Remark 3.2. To this point, we have developed an estimation theory based on state variable approximation only; that is, we used  $u^N$  to construct an approximate fit-to-data functional  $J^N$  which we then tacitly assumed could be minimized (numerically) to obtain an "optimal"  $\bar{\gamma}^N \in \Gamma$ . Of course, one cannot actually use a computer to implement such a parameter search since  $\Gamma$  is in fact a functional parameter set ( $\gamma \in \Gamma$  contains the functional components  $\phi_1, \phi_2, r, u_0, \kappa_1$ , and  $\kappa_2$ ). The problem of further approximating the parameter set  $\Gamma$  has been the subject of recent studies (see [10], [14], and [23]); we shall summarize, in particular, the results of [14] as they pertain to the problem at hand. For ease of presentation we shall assume that  $r, u_0, \kappa_1$ , and  $\kappa_2$  are known (there is an easy extension of these ideas to the case where these functional parameters are unknown) so that  $\gamma = (\xi, \phi_1, \phi_2)$  is the vector of parameters to be estimated. Since we use cubic B-spline approximations to approximate the functional parameters in our numerical examples (Section 4), we shall restrict our attention to a theory based on cubic splines only; a more general theory may be found in [14]. To this end, we take the (more regular) parameter set  $\Gamma$  to be a subset of

$$\tilde{\mathcal{S}}(m, \bar{m}) \equiv \{s = (\xi, \phi_1, \phi_2) \in [\delta, 1-\delta] \times C^1[0,1] \times$$

$$C^1[0,1] \mid 0 < m \leq \phi_i(x) \leq \bar{m} \text{ for } x \in [0,1],$$

$$\phi_i \in H^2(0,1), \text{ and } \|D^2 \phi_i\| \leq \bar{m}, i=1, 2\},$$

( $\delta \in (0,1)$  is fixed) and assume that  $\Gamma$  is compact in the  $R \times C^1[0,1] \times C^1[0,1]$  topology.

For each  $M$  we define the finite-dimensional (approximate) parameter sets  $\Gamma^M$  by  $\Gamma^M \equiv i^M(\Gamma)$ ; here  $i^M : R \times C^1[0,1] \times C^1[0,1] \rightarrow R \times C^2[0,1] \times C^2[0,1]$

is given by  $i^M(\xi, \phi_1, \phi_2) \equiv (\xi, \mathcal{I}^{M,\xi} \phi_1, \mathcal{I}^{M,\xi} \phi_2)$  where  $\mathcal{I}^{M,\xi} \phi$  is defined to be the (unique) cubic spline function  $\tilde{\phi}(x)$  satisfying  $\tilde{\phi}(x_k^M) = \phi(x_k^M)$ ,  $k = 0, 1, \dots, 2M$ , and  $\mathcal{D}\tilde{\phi}(x_0^M) = \mathcal{D}\phi(x_0^M)$ ,  $\mathcal{D}\tilde{\phi}(x_{2M}^M) = \mathcal{D}\phi(x_{2M}^M)$  (see, for example, Chapter 4 of [29]). The knots  $x_k^M$  are the  $\xi$ -dependent knots described earlier in this section, i.e.,  $x_k^M = k\xi/M$ ,  $k=0, \dots, M$ ,  $x_k^M = \xi + (k-M)(1-\xi)/M$  for  $k = M+1, \dots, 2M$ . We remark that we are also able to construct a parameter approximation scheme based on a uniform mesh (of mesh length  $1/M$ ) for components  $\phi_1^M, \phi_2^M$  of  $(\xi, \phi_1^M, \phi_2^M) \in \Gamma^M$ ; however, as is true with the approximation of state variables, the resulting numerical scheme is greatly simplified if the mesh depends on  $\xi$ , as well as on  $M$ . We shall defer to Section 4 a more detailed discussion on computational features of the resulting algorithm.

It is not difficult to use the ideas of [29; Chapter 4] to argue that, for fixed  $M$ , the mapping  $(\xi, \phi) \rightarrow \mathcal{I}^{M,\xi} \phi : [\delta, 1-\delta] \times C^1[0,1] \rightarrow C[0,1]$  is continuous and thus  $\Gamma^M$  is compact in the  $\mathbb{R} \times C[0,1] \times C[0,1]$  topology. In addition, given  $\varepsilon \in (0, m)$ , we may use a variation of [29; Theorem 4.5] to see that, for  $M > \mathfrak{m}$  sufficiently large (the choice of  $\mathfrak{m}$  is independent of  $\xi$ ),

$$0 < m - \varepsilon \leq \mathcal{I}^{M,\xi} \phi \leq \bar{m} + \varepsilon,$$

for all  $\xi \in [\delta, 1-\delta]$  and all  $\phi \in H^2(0,1)$ ,  $|\mathcal{D}^2 \phi| \leq \bar{m}$ ,  $m \leq \phi(x) \leq \bar{m}$ . Therefore, for  $M$  sufficiently large,  $\Gamma^M$  is a parameter set satisfying all conditions needed (namely  $\Gamma^M \subseteq \mathcal{S}(m-\varepsilon, \bar{m}+\varepsilon)$  and  $\Gamma^M$  satisfies hypothesis (H1)) order to apply the parameter estimation/state approximation theory developed thus far. In particular, for each  $N$  and  $M \geq \mathfrak{m}$ , there exists a solution  $\gamma^{N,M}$  to the problem of minimizing  $J^N$  over  $\Gamma^M$ . From the construction of  $\Gamma^M$  (and the compactness of  $\Gamma$ ) we know that there exists a sequence  $\{\hat{\gamma}^{N,M}\}$  in  $\Gamma$  with  $\bar{\gamma}^{N,M} = i^M(\hat{\gamma}^{N,M})$



and a subsequence  $\{\hat{\gamma}^{N_k, M^j}\}$  such that  $\hat{\gamma}^{N_k, M^j} \rightarrow \gamma^* \in \Gamma$  in the  $R \times C[0,1] \times C[0,1]$  topology. In addition,  $\bar{\gamma}^{N_k, M^j}$  satisfies

$$J^{N_k}(\bar{\gamma}^{N_k, M^j}) \leq J^{N_k}(\gamma), \quad \gamma \in \Gamma^{M^j}$$

from which it follows immediately that

$$(3.8) \quad J^{N_k}(\bar{\gamma}^{N_k, M^j}) \leq J^{N_k}(i^{M^j}(\hat{\gamma})), \quad \hat{\gamma} \in \Gamma.$$

Applying arguments very similar to those in [14], we may use the convergence of  $|i^M(\gamma) - \gamma|_\infty \rightarrow 0$  as  $M \rightarrow \infty$ , uniform in  $\gamma \in \Gamma$  [29; Theorem 4.5] to see that  $|\bar{\gamma}^{N_k, M^j} - \gamma^*|_\infty \rightarrow 0$  as  $N_k, M^j \rightarrow \infty$  and that, passing to the limit in (3.8),

$$(3.9) \quad J(\gamma^*) \leq J(\hat{\gamma}), \quad \hat{\gamma} \in \Gamma.$$

The parameter  $\gamma^*$  is thus a minimizer for  $J$  over  $\Gamma$ . We summarize these findings below.

Theorem 3.4 Let  $\Gamma^M \equiv i^M(\Gamma)$  and let  $\bar{\gamma}^{N, M}$  denote a solution to the problem of minimizing  $J^N$  over  $\Gamma^M$ . Then there is a subsequence  $\{\bar{\gamma}^{N_k, M^j}\}$  of  $\{\bar{\gamma}^{N, M}\}$  such that  $\bar{\gamma}^{N_k, M^j} \rightarrow \gamma^*$ , where  $\gamma^*$  is a solution to the problem of minimizing  $J$  over  $\Gamma$ . In fact, any convergent subsequence has as its limit a solution to the original estimation problem.

### 3.1. Approximate Estimation Problems Associated with "Discrete" Data

It is possible, under additional smoothness assumptions on solutions, to use variational-type estimates (similar to those found above or in [10], [30]) to argue pointwise (in  $x$ ) convergence of state variables; i.e.,  $u^N(t, x; \gamma^N) \rightarrow u(t, x; \bar{\gamma})$  whenever  $\gamma^N \rightarrow \bar{\gamma}$ , for  $(t, x) \in (0, T) \times [0, 1]$ . Results of this type

lead naturally to a statement about the approximation of a solution for the problem of minimizing the "pointwise" fit-to-data criterion  $\tilde{J}$  (see (2.3)) over  $\Gamma$ . We shall briefly summarize our findings below.

It is not surprising that for this formulation we require additional smoothness assumptions on the parameters. Specifically we take  $\gamma = (s, u_0) \in \Gamma$  where  $\Gamma$  is a subset of  $\hat{\Gamma} = \hat{\mathcal{S}} \times H_0^1(0,1)$ , and  $\hat{\mathcal{S}} \equiv \{s = (\xi, \phi_1, \phi_2, r) \in \mathcal{S} \mid |\mathcal{D}\phi_i|_\infty \leq \bar{m}, i = 1,2\}$ .

In addition to hypotheses (H1) - (H3), we make the following assumption (which in general may impose additional constraints on parameters and the applied force  $f$ ):

(H4) For any  $\gamma \in \Gamma$ , the mapping  $s \rightarrow u_t(s; \gamma)$  is in  $L_2((0,T); H_0^1(0,1))$ .

Defining approximate state spaces  $X^N(q)$  as before, we seek approximations  $u^N$  to  $u$ , where, for any given  $\gamma = (\xi, \phi_1, \phi_2, r, u_0) \in \Gamma$ ,  $q = \phi_1 + H_\xi \phi_2$ ,

$u^N(t) = u^N(t; \gamma)$  satisfies

$$(3.10) \quad \begin{cases} \langle u_t^N(t), v \rangle = - \langle q \mathcal{D}u^N(t), \mathcal{D}v \rangle + \langle F(t; r), v \rangle, & t \in (0, T), \\ u^N(0) = \mathcal{P}^N(q)u_0 \end{cases}$$

for all  $v \in X^N(q)$ ;  $\mathcal{P}^N$  differs from  $P^N$  defined in (3.2) in that  $\mathcal{P}^N: H_0^1(0,1) \rightarrow X^N(q)$  is the orthogonal projection in the  $H_0^1(0,1)$  (rather than  $L_2(0,1)$ ) topology. It is not difficult to see that there exists a unique solution  $u^N(\gamma)$  of (3.10). In addition, for each  $(t, x) \in [0, T] \times [0, 1]$ , the mapping  $\gamma \rightarrow u^N(t; \gamma)(x)$  is continuous (in the  $\mathcal{S} \times H_0^1(0,1)$  topology on  $\gamma$ ).

The convergence result that follows is a pointwise analog of Lemma 3.2.

Lemma 3.3. Let  $\{\gamma^N\}$  be given in  $\Gamma$  such that  $\gamma^N \rightarrow \bar{\gamma}$  in the  $\mathcal{S} \times H_0^1(0,1)$  topology,  $\bar{\gamma} = (\bar{\xi}, \bar{\phi}_1, \bar{\phi}_2, \bar{r}, \bar{u}_0) \in \Gamma$ . If  $\bar{u}_0 \in V_{\bar{q}}$ , then  $u^N(t; \gamma^N) \rightarrow u(t; \bar{\gamma})$  (as  $N \rightarrow \infty$ ) in the  $H_0^1(0,1)$  topology, uniform in  $t \in [0, T]$ ; here  $u^N(\gamma^N)$  and  $u(\bar{\gamma})$  are solutions of (3.10) and (2.7) corresponding to  $\gamma^N$  and  $\bar{\gamma}$ , respectively.

Proof: We shall write  $u^N = u^N(t; \gamma^N)$ ,  $u = u(t; \bar{\gamma})$ ,  $q^N = \phi_1^N + H_{\xi}^N \phi_2^N$ ,  $\bar{q} = \bar{\phi}_1 + H_{\bar{\xi}} \bar{\phi}_2$ ,  $F^N = F(t; r^N)$ , and  $F = F(t; \bar{r})$  throughout. We note that

$$|\mathcal{D}(u^N - u)| \leq |\mathcal{D}(u^N - \mathcal{I}^N u)| + |\mathcal{D}(\mathcal{I}^N u - u)|$$

so that, using (3.6) and the fact that  $u(t) \in V_{\bar{q}}$  for  $t \geq 0$ , it suffices to show  $|\mathcal{D}(u^N - \mathcal{I}^N u)| \rightarrow 0$  as  $N \rightarrow \infty$ .

Using (3.10), (2.7), and arguments similar to those in the proof of Lemma 3.2, we may argue that, for  $v \in X^N = X^N(q^N)$ ,

$$\begin{aligned} \langle u_t^N - \mathcal{I}^N u_t, v \rangle + \langle q^N \mathcal{D} u^N - q^N \mathcal{D} \mathcal{I}^N u, \mathcal{D} v \rangle &= \langle u_t - \mathcal{I}^N u_t, v \rangle + \langle \bar{q} \mathcal{D} u - q^N \mathcal{D} \mathcal{I}^N u, \mathcal{D} v \rangle \\ &\quad + \langle F^N - F, v \rangle \end{aligned}$$

and in particular, using  $v = u_t^N - \mathcal{I}^N u_t$ ,

$$\begin{aligned} |u_t^N - \mathcal{I}^N u_t|^2 + \langle q^N \mathcal{D}(u^N - \mathcal{I}^N u), \mathcal{D}(u_t^N - \mathcal{I}^N u_t) \rangle \\ \leq \frac{1}{2} |u_t - \mathcal{I}^N u_t|^2 + |u_t^N - \mathcal{I}^N u_t|^2 \\ + \frac{1}{2} |F^N - F|^2 + \langle \bar{q} \mathcal{D} u - q^N \mathcal{D} \mathcal{I}^N u, \mathcal{D}(u_t^N - \mathcal{I}^N u_t) \rangle. \end{aligned}$$

We thus find that

$$\begin{aligned}
& \frac{1}{2} \frac{d}{dt} \{ |(q^N)^{1/2} \mathfrak{D}(u^N - \mathcal{I}^N u)|^2 - 2 \langle \bar{q} \mathfrak{D} u - q^N \mathfrak{D} \mathcal{I}^N u, \mathfrak{D}(u^N - \mathcal{I}^N u) \rangle \} \\
& \leq \frac{1}{2} |u_t - \mathcal{I}^N u_t|^2 + \frac{1}{2} |F^N - F|^2 - \langle \bar{q} \mathfrak{D} u_t - q^N \mathfrak{D} \mathcal{I}^N u_t, \mathfrak{D}(u^N - \mathcal{I}^N u) \rangle
\end{aligned}$$

where we have used the fact that  $\frac{d}{dt} \mathcal{I}^N u = \mathcal{I}^N u_t$ ; it therefore follows that

$$\begin{aligned}
(3.11) \quad m |\mathfrak{D}(u^N - \mathcal{I}^N u)(t)|^2 & \leq \frac{2}{m} \sigma_1^N(t) + \frac{m}{2} \sigma_2^N(t) + \sigma_1^N(0) + (\bar{m}+1) \sigma_2^N(0) + \sigma_3^N \\
& + \sigma_4^N + \sigma_5^N + \sigma_6^N,
\end{aligned}$$

where

$$\sigma_1^N(t) = |\bar{q} \mathfrak{D} u(t) - q^N \mathfrak{D} \mathcal{I}^N u(t)|^2,$$

$$\sigma_2^N(t) = |\mathfrak{D}(u^N - \mathcal{I}^N u)(t)|^2,$$

$$\sigma_3^N = \int_0^T |u_t - \mathcal{I}^N u_t|^2,$$

$$\sigma_4^N = \int_0^T |F^N - F|^2$$

$$\sigma_5^N = -2 \int_0^T \langle (\bar{q} - q^N) \mathfrak{D} u_t, \mathfrak{D}(u^N - \mathcal{I}^N u) \rangle,$$

and

$$\sigma_6^N = -2 \int_0^T \langle q^N \mathfrak{D}(u_t - \mathcal{I}^N u_t), \mathfrak{D}(u^N - \mathcal{I}^N u) \rangle.$$

Using (3.6) and arguments like those for  $\tau_3^N$  in the proof of Lemma 3.2, we find that

$$\begin{aligned}\sigma_1^N(t) &\leq 2|(\bar{q} - q^N)\mathfrak{D}u(t)|^2 + 2\bar{m}^2|\mathfrak{D}(u - \mathcal{I}^N u)(t)|^2 \\ &\leq 2k(N)|\mathfrak{D}u(t)|^2 + 2\bar{m}^2(m\pi N)^{-2}|A(\bar{q})u(t)|^2,\end{aligned}$$

where

$$k(N) = 4(|\bar{\phi}_1 - \phi_1^N|_\infty^2 + |\bar{\phi}_2 - \phi_2^N|_\infty^2) + 2\bar{m}^2|\bar{\xi} - \xi^N| \rightarrow 0 \quad \text{as } N \rightarrow \infty.$$

Thus, for  $t \in [0, T)$ ,

$$\begin{aligned}\sigma_1^N(t) &\leq 2k(N) \sup_{t \in [0, T)} |\mathfrak{D}u(t)|^2 + 2\bar{m}^2(m\pi N)^{-2} \sup_{t \in [0, T)} |A(\bar{q})u(t)|^2 \\ &\equiv \tilde{\sigma}_1^N.\end{aligned}$$

Considering  $\sigma_2^N$  and  $\sigma_5^N$ , we find

$$\sigma_2^N(t) \leq \sup_{t \in [0, T)} |\mathfrak{D}(u^N - \mathcal{I}^N u)(t)|^2$$

and

$$\begin{aligned}\sigma_5^N &\leq \frac{4T}{m} \int_0^T |(\bar{q} - q^N)\mathfrak{D}u_t|^2 + \frac{m}{4T} \int_0^T |\mathfrak{D}(u^N - \mathcal{I}^N u)|^2 \\ &\leq \frac{4T}{m} k(N) \int_0^T |\mathfrak{D}u_t|^2 + \frac{m}{4} \sup_{t \in [0, T)} |\mathfrak{D}(u^N - \mathcal{I}^N u)(t)|^2.\end{aligned}$$

In addition, we may integrate by parts to show that

$$\begin{aligned}\sigma_6^N &= -2 \sum_{k=1}^{2N} \int_0^T \int_{x_{k-1}^N}^{x_k^N} \mathfrak{D}(u_t - \mathcal{I}^N u_t) q^N \mathfrak{D}(u^N - \mathcal{I}^N u) \\ &= 2 \sum_{k=1}^{2N} \int_0^T \int_{x_{k-1}^N}^{x_k^N} (u_t - \mathcal{I}^N u_t) \{ \mathfrak{D}q^N \mathfrak{D}(u^N - \mathcal{I}^N u) + q^N \cdot 0 \}\end{aligned}$$

$$\begin{aligned}
&\leq \bar{m} \left\{ 8\bar{m}Tm^{-1} \int_0^T |u_t - \mathcal{I}^N u_t|^2 + \frac{m}{8\bar{m}T} \int_0^T |\mathcal{D}(u^N - \mathcal{I}^N u)|^2 \right\} \\
&\leq 8\bar{m}^2 Tm^{-1} \int_0^T |u_t - \mathcal{I}^N u_t|^2 + \frac{m}{8} \sup_{t \in [0, T]} |\mathcal{D}(u^N - \mathcal{I}^N u)(t)|^2,
\end{aligned}$$

where we have used  $|\mathcal{D}\phi_i^N| \leq \bar{m}$ ,  $i = 1, 2$ . Therefore, from (3.11),

$$\begin{aligned}
\frac{m}{8} \sup_{t \in [0, T]} |\mathcal{D}(u^N - \mathcal{I}^N u)(t)|^2 &\leq (2m^{-1} + 1)\tilde{\sigma}_1^N + (\bar{m}+1)|\mathcal{D}(\mathcal{P}^N u_0^N - \mathcal{I}^N \bar{u}_0)|^2 \\
&\quad + (8\bar{m}^2 Tm^{-1} + 1)\sigma_3^N + \sigma_4^N \\
&\quad + 4Tk(N)m^{-1} \int_0^T |\mathcal{D}u_t|^2,
\end{aligned}$$

where  $\tilde{\sigma}_1^N, \sigma_4^N \rightarrow 0$  as  $N \rightarrow \infty$ . We may apply (H4) and standard spline estimates (see, for example, Theorem 2.4 of [29]) to also show that  $\sigma_3^N \rightarrow 0$  as  $N \rightarrow \infty$ . It remains to consider  $|\mathcal{D}(\mathcal{P}^N u_0^N - \mathcal{I}^N \bar{u}_0)|^2$ : We note, using the properties of  $\mathcal{P}^N$ , that

$$\begin{aligned}
|\mathcal{D}(\mathcal{P}^N u_0^N - \mathcal{I}^N \bar{u}_0)| &\leq |\mathcal{D}(\mathcal{P}^N u_0^N - \mathcal{P}^N \bar{u}_0)| + |\mathcal{D}(\mathcal{P}^N \bar{u}_0 - \bar{u}_0)| + |\mathcal{D}(\bar{u}_0 - \mathcal{I}^N \bar{u}_0)| \\
&\leq |\mathcal{D}(u_0^N - \bar{u}_0)| + 2|\mathcal{D}(\mathcal{I}^N \bar{u}_0 - \bar{u}_0)|,
\end{aligned}$$

where each term converges to 0 as  $N \rightarrow \infty$  from the convergence of  $u_0^N$  to  $\bar{u}_0$  in  $H_0^1$  and (3.6). The proof of the lemma thus obtains.

Finally, we prove a more general convergence result.

Theorem 3.5. Let  $\{\gamma^N\}$  be given in  $\Gamma$  with  $\gamma^N \rightarrow \bar{\gamma} \in \Gamma$  in the  $\mathcal{S} \times H_0^1(0,1)$  topology. Then, for each  $t \in [0,T)$ ,

$$u^N(t; \gamma^N) \rightarrow u(t; \bar{\gamma}) \text{ in } H_0^1(0,1)$$

as  $N \rightarrow \infty$ .

Proof: We first demonstrate the continuity of the mapping

$u_0 \rightarrow u(t; (s, u_0)) : H_0^1 \rightarrow H_0^1$ , uniform in  $t \in [0, T)$  and  $s \in \mathcal{S}$ . Let  $u_0, u'_0 \in H_0^1$  and  $u = u(t; (s, u_0))$ ,  $u' = u(t; (s, u'_0))$ . For any  $v \in H_0^1$ ,

$$\langle u_t - u'_t, v \rangle = - \langle q(Du - Du'), Dv \rangle$$

so that, using  $v = u_t - u'_t \in H_0^1$ ,

$$0 \leq |u_t - u'_t|^2 = - \frac{1}{2} \frac{d}{dt} |q^{1/2} D(u - u')|^2.$$

Applying the Gronwall inequality, we find that

$$|q^{1/2} D(u - u')(t)|^2 \leq |q^{1/2} D(u_0 - u'_0)|^2,$$

or that

$$|D(u - u')(t)|^2 \leq \bar{m}m^{-1} |D(u_0 - u'_0)|^2,$$

so that the continuous dependence result obtains.

We may construct similar arguments to demonstrate (using (3.10)) that

$$\begin{aligned} |D(u^N(t; u_0) - u^N(t; u'_0))|^2 &\leq \bar{m}m^{-1} |D(p^N u_0 - p^N u'_0)|^2 \\ &\leq \bar{m}m^{-1} |D(u_0 - u'_0)|^2 \end{aligned}$$

so that the mappings  $u_0 \rightarrow u^N(t; (s, u_0)) : H_0^1 \rightarrow H_0^1$  are also continuous, uniform in  $t$ ,  $s$ , and  $N$ .

Let  $\{\gamma^N\}$  be given in  $\Gamma$  with  $\gamma^N \rightarrow \bar{\gamma}$ , where  $\bar{\gamma} = (\bar{\xi}, \bar{\phi}_1, \bar{\phi}_2, \bar{r}, \bar{u}_0) \in \Gamma$ . To argue the convergence of  $u^N(t; \gamma^N) \rightarrow u(t; \bar{\gamma})$  in  $H_0^1$  for arbitrary  $\bar{u}_0 \in H_0^1$  (using estimates like those in the proof of Theorem 3.2) we need only demonstrate that  $V_{\bar{q}}$ ,  $\bar{q} = \bar{\phi}_1 + H_{\bar{\xi}} \bar{\phi}_2$ , is dense in  $H_0^1$ . To this end, let  $\varepsilon > 0$  and  $\psi$  be arbitrary in  $H_0^1(0,1)$  and define  $\tilde{\psi} = \psi - h$  where  $h$  is given by

$$h(x) = \begin{cases} \psi(\bar{\xi})x/\bar{\xi}, & x \in [0, \bar{\xi}] \\ p(x), & x \in (\bar{\xi}, 1] \end{cases}.$$

Here  $p$  is a quadratic polynomial satisfying  $p(1) = 0$ ,  $p(\bar{\xi}) = \psi(\bar{\xi})$ , and  $\bar{q}p(\bar{\xi}) = \bar{q}(\bar{\xi}^-)\psi(\bar{\xi})/\bar{\xi}\bar{q}(\bar{\xi}^+)$  (it is easy to show such a  $p$  exists, under the condition  $\bar{\xi} \neq 1$ ). It is clear from the construction of  $h$  that we have  $h \in V_{\bar{q}}$ ,  $\tilde{\psi} \in H_0^1(0,1)$ , and  $\tilde{\psi}(x) = 0$  for  $x = 0, \bar{\xi}, 1$ , so that  $\tilde{\psi} \in H_0^1(0, \bar{\xi})$ ,  $\tilde{\psi} \in H_0^1(\bar{\xi}, 1)$ . From the definition of  $H_0^1$ , there exists  $\tilde{\zeta}_1 \in C_0^\infty(0, \bar{\xi})$  and  $\tilde{\zeta}_2 \in C_0^\infty(\bar{\xi}, 1)$  such that

$$|\tilde{\psi} - \tilde{\zeta}_1|_{H_0^1(0, \bar{\xi})}^2 < \varepsilon/2 \text{ and } |\tilde{\psi} - \tilde{\zeta}_2|_{H_0^1(\bar{\xi}, 1)}^2 < \varepsilon/2. \text{ Finally, defining } \zeta \text{ on } [0, 1]$$

by  $\zeta = \tilde{\zeta}_1 + h$  on  $[0, \bar{\xi}]$ ,  $\zeta = \tilde{\zeta}_2 + h$  on  $(\bar{\xi}, 1]$ , it is easy to see that  $\zeta \in V_{\bar{q}}$

( $\zeta \in H_0^1(0,1)$ ,  $\bar{q}D\zeta \in H^1(0, \bar{\xi})$ ,  $\bar{q}D\zeta \in H^1(\bar{\xi}, 1)$ , and  $\bar{q}D\zeta$  is continuous at  $x = \bar{\xi}$ )

and  $|\psi - \zeta|_{H_0^1(0,1)}^2 = |\tilde{\psi} - \tilde{\zeta}_1|_{H_0^1(0, \bar{\xi})}^2 + |\tilde{\psi} - \tilde{\zeta}_2|_{H_0^1(\bar{\xi}, 1)}^2 < \varepsilon$ . Therefore,  $V_{\bar{q}}$  is

dense in  $H_0^1(0,1)$  and the proof of the theorem is complete.

Finally, we return to the problem of (approximately) determining a minimizer  $\tilde{\gamma}^*$  for the (pointwise) least squares criterion  $\tilde{J}$ . The proof of our final result uses the estimates derived in this section, following the proof of Theorem 3.3.



Theorem 3.6. For each  $N$ , let  $\tilde{\gamma}^N$  denote a solution for the problem of minimizing  $\tilde{J}^N$  over  $\Gamma$ , where  $\tilde{J}^N(\gamma) = \sum_{i=1}^n \sum_{j=1}^{\tilde{n}} |C(t_i, x_j; \gamma) u^N(t_i, x_j; \gamma) - \hat{u}_{ij}|^2$  and  $u^N$  is the solution of (3.10) associated with  $\gamma \in \Gamma$ . Then there exists  $\tilde{\gamma}^* \in \Gamma$  and a subsequence  $\{\tilde{\gamma}^{N_k}\}$  of  $\{\tilde{\gamma}^N\}$  such that  $\tilde{\gamma}^{N_k} \rightarrow \tilde{\gamma}^*$ ,  $\tilde{J}^{N_k}(\tilde{\gamma}^{N_k}) \rightarrow \tilde{J}(\tilde{\gamma}^*)$ , and  $\tilde{\gamma}^*$  is a solution for the problem of minimizing  $\tilde{J}$  over  $\Gamma$ .

One may also easily modify the arguments given above to include  $\rho$  as a parameter (see Remark 3.1) and to prove a "double approximation" result similar to Theorem 3.4 for both the state and estimated (functional) parameters.

#### 4. Implementation and Numerical Findings

A desirable feature of the spline-based scheme developed in preceding sections is the ease of implementation of the approximation ideas, especially when the points of discontinuity  $\xi_i$ ,  $i=1, \dots, \mu-1$ , for coefficients are unknown and to be estimated. In what follows we describe how the particular state approximation framework chosen here serves to facilitate (from a computational standpoint) the parameter estimation/approximation process. We conclude the section by presenting our findings for some representative test examples.

We begin by examining the approximating ordinary differential equation (3.1) rewritten here in terms of  $w^N(t; \gamma) \equiv (w_1^N(t; \gamma), w_2^N(t; \gamma), \dots, w_{2N-1}^N(t; \gamma))^T$ , where the  $w_i^N$ , defined in Section 3, are the coefficients in the expansion  $u^N(t; \gamma) = \sum_{i=1}^{2N-1} w_i^N(t; \gamma) B_i^N(q)$ . Using this notation, the ODE may be written

$$(4.1) \quad \begin{cases} Q^N \dot{w}^N(t) = -K^N w^N(t) + G^N(t), & t \in (0, T), \\ w^N(0) = w_0^N; \end{cases}$$

here the  $(2N-1)$ -square matrices  $Q^N = Q^N(\gamma)$  and  $K^N = K^N(\gamma)$  have entries

$$Q_{i,j}^N = \langle B_j^N(q), B_i^N(q) \rangle, \quad$$

$$K_{i,j}^N = \langle q \mathcal{D} B_j^N(q), \mathcal{D} B_i^N(q) \rangle, \quad$$

while the perturbation term and initial condition satisfy

$$G^N(t) = G^N(t; \gamma) \equiv (\langle F(t; r), B_1^N(q) \rangle, \dots, \langle F(t; r), B_{2N-1}^N(q) \rangle)^T$$

and

$$w_0^N = w_0^N(\gamma) \equiv (Q^N)^{-1} (\langle u_0, B_1^N(q) \rangle, \dots, \langle u_0, B_{2N-1}^N(q) \rangle)^T,$$

respectively.

Implementation of a computational scheme to estimate (approximate) parameters in (4.1) is greatly facilitated by the choice of basis elements for the  $N^{\text{th}}$  approximate state space  $X^N(q)$ . In order to best indicate some of the advantages of this approximation framework, we first consider the special case where only  $q$  is unknown, where  $q = \phi_1 + H_\xi \phi_2$ , and  $\phi_1$  and  $\phi_2$  are constants. First, it is easy to see how our choice of a linear spline approximation scheme yields matrices  $K^N$  and  $Q^N$  that are quite simple in structure: For a given value of  $q$ , the inner products appearing in these matrices may be determined from explicit formulas (depending on  $N$  and  $\xi$ ), a few of which are given here. For example, diagonal entries in the (tridiagonal) matrices  $Q^N$  and  $K^N$  are given by

$$Q_{i,i}^N = 2\xi/3N, \quad i = 1, \dots, N-1,$$

$$Q_{N,N}^N = 1/3N,$$

$$Q_{i,i}^N = 2(1-\xi)/3N, \quad i = N+1, \dots, 2N-1,$$

$$(4.2) \quad K_{i,i}^N = 2N\phi_1/\xi, \quad i = 1, \dots, N-1,$$

$$(4.3) \quad K_{N,N}^N = N\phi_1/\xi + N(\phi_1 + \phi_2)/(1-\xi),$$

and

$$(4.4) \quad K_{i,i}^N = 2N(\phi_1 + \phi_2)/(1-\xi), \quad i = N+1, \dots, 2N-1,$$

with similar representations for off-diagonal elements. We note that we are able to avoid time-consuming and error-producing numerical quadratures; in addition, our approach is more desirable (from a computational point of view) than a method based on a uniform mesh size. For example, if for each  $N$  we simply subdivide  $[0,1]$  into units of length  $\frac{1}{N}$  (so that position of  $\xi$  is not taken into account) the matrix  $Q^N$  will be fixed throughout the estimation process; this however is at the expense of considerable added difficulties associated with

evaluating entries in  $K^N$ . Using a uniform mesh, some of the inner products must be "broken up" at the point  $\xi$ , e.g.,

$$\langle q \mathcal{D}B_j^N, \mathcal{D}B_i^N \rangle = \phi_1 \int_0^\xi \mathcal{D}B_j^N \mathcal{D}B_i^N + (\phi_1 + \phi_2) \int_\xi^1 \mathcal{D}B_j^N \mathcal{D}B_i^N,$$

requiring (multiple) numerical quadratures every time that  $q$  (and thus  $\xi$ ) is updated. In contrast, with the  $\xi$ -dependent structure chosen here we need only recombine simple algebraic expressions (such as those given in (4.2) - (4.4)) to obtain the elements of  $K^N$ .

Many of these computational advantages are still present in the case where the  $\phi_i$  are not assumed to be constant. If, for example,  $M$  and  $N$  are fixed and  $\Gamma^M$  consists of cubic spline element approximations for  $\phi_1, \phi_2$  (defined  $\xi$ -dependent mesh of points  $x_k^M$ ; see Section 3), many of the quadratures may still be performed in advance of the iterative process. In particular, if we let  $\phi_n^M(x) = \sum_{m=1}^{k(M)} \gamma_{n,m}^M \mathcal{G}_m^M(x)$ , for  $n=1,2$ , where  $\mathcal{G}_m^M$  are the usual cubic B-spline basis elements defined using the mesh points  $\{x_k^M, k=0, \dots, 2M\}$ , we find that  $K_{i,j}^N = \langle q \mathcal{D}B_j^N, \mathcal{D}B_i^N \rangle$  may now be written as

$$(4.5) \quad K_{i,j}^N = \sum_{m=1}^{k(M)} \gamma_{1,m}^M \langle \mathcal{G}_m^M \mathcal{D}B_j^N, \mathcal{D}B_i^N \rangle_{L_2(0,\xi)} + \sum_{m=1}^{k(M)} (\gamma_{1,m}^M + \gamma_{2,m}^M) \langle \mathcal{G}_m^M \mathcal{D}B_j^N, \mathcal{D}B_i^N \rangle_{L_2(\xi,1)}.$$

Since simple explicit algebraic expressions (in terms of  $\xi$  and  $M$ ) exist for  $\mathcal{G}_m^M$ , the quadratures in (4.5) may also be worked out easily in advance (yielding expressions involving  $\xi$ ,  $M$ , and  $N$ ); as  $q$ -iterates are updated (i.e., coefficients  $\gamma_{n,m}^M$  are updated) it becomes a simple task to calculate the new entries in  $K^N$ .

We consider here numerical examples where  $\gamma^*$  is known and we have generated synthetic data for use in testing our ideas. In all examples presented here, we assume that  $r$  and  $u_0$  are known and fixed at their true values so that only  $q = \phi_1 + H_\xi \phi_2$  is unknown (i.e.,  $\gamma = (\xi, \phi_1, \phi_2)$ ) and to be determined. The special

problems associated with estimating this discontinuous coefficient have been the focus of our efforts throughout; the problem of identifying continuous functional parameters and initial conditions has been considered elsewhere [13], [14], [16]. For each example that follows, both  $\gamma^*$  and  $u(\gamma^*)$  are selected in advance while the appropriate forcing function  $f$  is artificially determined by substituting  $\gamma^*$ ,  $u(\gamma^*)$  into (2.1). For chosen sample times  $t_i$ ,  $i=1, \dots, n$ , and sampling locations  $x_j$ ,  $j=1, \dots, \tilde{n}$  (discrete data is used for these examples), data is generated by setting  $\hat{u}_{ij} = u(t_i, x_j; \gamma^*)$ , with random noise added in some cases. We note that the sample data is not generated using our spline-based scheme; rather, the data is constructed from an analytic expression for the solution and thus is independent of the methods we illustrate here.

We begin the parameter estimation process by supplying an initial guess of  $\gamma^0$  to IMSL's minimization routine ZXSSQ (a Levenberg-Marquardt algorithm) which numerically attempts to determine a minimum, for given  $N$ , to  $\tilde{J}^N$  (using  $\epsilon \equiv 1$  in (2.3)) over a fixed constraint set  $\Gamma^M$ . Here  $u^N(\gamma)$  is the solution to (3.1) calculated using IMSL's DGEAR, an ODE solver, where the known values of  $u_0$  and  $f$  are used in the equations. We note that although we are actually using the cost functional associated with discrete observations  $\hat{u}_{ij}$ , the approximating equations (4.1) differ somewhat from those defined in Section 3.1: Indeed it is not surprising that, in practice, we obtain pointwise convergence of the approximating states under hypotheses more general than those needed in Section 3.1 so that we may, in fact, relax some of the restrictions on the approximating system.

Example 4.1. In our first example we take

$$q^*(x) = \begin{cases} 15. & , \quad 0 \leq x < .6 \\ 50. & , \quad .6 \leq x \leq 1 \end{cases}$$

and define  $u(t, \cdot; \gamma^*) \in \text{dom } A(q^*)$  by

$$u(t, x; \gamma^*) = \begin{cases} x(70-100x)(t^2+2) & , \quad 0 \leq x < .6 \\ (15-15x)(t^2+2) & , \quad .6 \leq x \leq 1. \end{cases}$$

In examples 4.1.a - 4.1.c below we seek to estimate  $\gamma = (\xi, \phi_1, \phi_2) \in \Gamma \subseteq R^3$  (with true value  $\gamma^* = (.6, 15., 50.)$ ) using an initial guess of  $\gamma^0 = (.8, 30., 30.)$ .

In each case we obtain the converged values  $\bar{\gamma}^N$  for  $N=4, 8, 16$ , and  $24$ , using  $\gamma^0$  to start the iterative scheme for  $N=4$ , and previous converged values as start-ups for  $N=8, 16, 24$  (e.g.,  $\bar{\gamma}^4$  is used as initial guess for the  $N=8$  run).

Example 4.1.a. Data is generated for this example using  $\hat{u}_{ij} = u(t_i, x_j; \gamma^*)$  for  $t_i = .5i, i=1, \dots, 4$ , and  $x_j = .1j, j=1, \dots, 9$ . Our findings are reported in Table 4.1.a.

Example 4.1.b. We repeat the last example except that spatial sampling locations are now given by  $x_j = .1j + .05, j=0, 1, \dots, 9$  (so that there is no spatial observation point at  $\xi^*$ , the point of discontinuity). We summarize our results in Table 4.1.b. and note that there is little change between this example and Example 4.1.a.

Example 4.1.c. We repeat Example 4.1.a, but add noise to the data. In this case we define  $\hat{u}_{ij} = u(t_i, x_j; \gamma^*) + r_{ij}$  where  $\{r_{ij}\}$  are Gaussian random numbers which (with 98% certainty) fall in the range  $[-.06\bar{u}, .06\bar{u}]$ ,  $\bar{u} = \sum_{i,j} \hat{u}_{ij} / (n\tilde{n})$ . Our findings for this example are summarized in Table 4.1.c.

In the examples that follow we shall shorten our discussion by abbreviating the length (and number) of tables and by displaying some results graphically. The rather detailed presentation given for Example 4.1 was provided simply for the purpose of observing if noise in the data or changes in the placement of data affected the outcome.

Table 4.1.a. -- Example 4.1.a.

<u>N</u>	<u><math>\bar{\xi}^N</math></u>	<u><math>\bar{\phi}_1^N</math></u>	<u><math>\bar{\phi}_2^N</math></u>	<u><math>\tilde{J}^N</math></u>	<u>CP time (secs)</u>	<u>No. of iterates</u>
4	.623	14.669	51.950	$1.5 \times 10^2$	28.	13
8	.602	14.845	50.672	$1.5 \times 10^0$	54.	7
16	.600	14.961	50.095	$8.8 \times 10^{-2}$	202.	7
24	.600	15.000	50.000	$5.6 \times 10^{-9}$	141.	4

Table 4.1.b: -- Example 4.1.b.

<u>N</u>	<u><math>\bar{\xi}^N</math></u>	<u><math>\bar{\phi}_1^N</math></u>	<u><math>\bar{\phi}_2^N</math></u>	<u><math>\tilde{J}^N</math></u>	<u>CP time (secs)</u>	<u>No. of iterates</u>
4	.621	14.956	48.494	$9.0 \times 10^1$	32.	20
8	.607	15.009	50.063	$3.8 \times 10^1$	35.	7
16	.601	14.991	49.728	$1.2 \times 10^{-1}$	355.	13
24	.600	15.000	50.000	$5.7 \times 10^{-9}$	239.	5

Table 4.1.c. -- Example 4.1.c. (Noisy data)

<u>N</u>	<u><math>\bar{\xi}^N</math></u>	<u><math>\bar{\phi}_1^N</math></u>	<u><math>\bar{\phi}_2^N</math></u>	<u><math>\tilde{J}^N</math></u>	<u>CP time (secs)</u>	<u>No. of iterates</u>
4	.621	14.730	51.573	$1.6 \times 10^2$	27.	10
8	.599	14.887	50.434	$8.1 \times 10^0$	68.	10
16	.598	14.991	50.296	$8.0 \times 10^0$	178.	5
24	.597	15.006	50.149	$8.3 \times 10^0$	733.	8

In Example 4.2 below, we illustrate the use of our methods in problems with two discontinuities  $\xi_1, \xi_2$ , in  $q$ ; the example also serves to illustrate that we are able to accurately estimate  $\xi_i$  even when the forcing function  $f$  does not contain discontinuities at each of those points.

Example 4.2. We seek here the "true" value of  $q$  given by

$$q^* = \begin{cases} 1.0 & , \quad 0. \leq x < .2 \\ 6.0 & , \quad .2 \leq x < .6 \\ 0.5 & , \quad .6 \leq x \leq 1. \end{cases}$$

In this case,  $\gamma^* \equiv (\xi_1^*, \xi_2^*, \phi_1^*, \phi_2^*, \phi_3^*) = (.2, .6, 1., 6., .5)$  and the true solution corresponding to  $\gamma^*$  is

$$u(t, x; \gamma^*) = \begin{cases} 30x & 0 \leq x < .2 \\ 5x + 5 & .2 \leq x < .6 \\ -200x^2 + 300x - 100 & .6 \leq x \leq 1 \end{cases},$$

with data available at  $t_i = .5i$ ,  $i=1, \dots, 4$ , and  $x_j = .1j$ ,  $j=1, \dots, 9$ .

A sample of our findings is given in Table 4.2 below, where the converged values reported were obtained after 501 CP seconds, with  $\tilde{J}^8 = 8.3 \times 10^{-6}$ .

Table 4.2 -- Example 4.2

	<u><math>\xi_1</math></u>	<u><math>\xi_2</math></u>	<u><math>\phi_1</math></u>	<u><math>\phi_2</math></u>	<u><math>\phi_3</math></u>
<u>Initial guess:.</u>	.300	.700	5.000	5.000	5.000
<u>Converged values</u> <u>(N = 8):</u>	.200	.600	1.000	6.000	.5000
<u>True values:</u>	.200	.600	1.000	6.000	.5000



We consider now two examples where the "true"  $q^* = \phi_1^* + H_{\xi} \phi_2^*$  involves nonconstant values of  $\phi_1^*$  and  $\phi_2^*$ . In each case we search for approximate  $\phi_1$  and  $\phi_2$  in the cubic spline space constructed using an  $M=1$  level of approximation (see Section 3).

Example 4.3. Here we seek to estimate the "true" parameter

$$q^* = \begin{cases} 2x + 12 & , \quad 0 \leq x < .6 \\ 1100x^2/9 & , \quad .6 \leq x < 1 \end{cases} ,$$

starting from the initial guess for  $q$  of  $q^0 \equiv 3$  on  $[0,1]$  (with start-up value for  $\xi$  of  $\xi = .5$ ). The solution

$$u(t, x; \gamma^*) = \begin{cases} (70x - 100x^2)(t^2 + 2) & , \quad 0 \leq x < .6 \\ 15(1-x)(t^2 + 2) & , \quad .6 \leq x \leq 1 \end{cases} ,$$

is used to generate data at  $t_i = .5i$ ,  $i=1, \dots, 4$ , and  $x_j = .1j$ ,  $j=1, \dots, 9$ .

In Figure 1a we compare the estimated  $\bar{q}^{N,M} = \bar{\phi}_1^{N,M} + H_{\xi N} \bar{\phi}_2^{N,M}$  ( $N=16$ ,  $M=1$ ) with the "true" coefficient  $q^*$ . Figure 1b is the same graph that has been enlarged and restricted to the interval  $[.4, .63]$  in order to better distinguish between "true" and approximate curves.

Example 4.4. Again we estimate a functional parameter with true representation given by

$$q^* = \begin{cases} 27.424 - 40x & , \quad 0 \leq x < .3 \\ 90(x-.3)^2 + 18 & , \quad .3 \leq x \leq 1. \end{cases}$$

Data is generated as in Example 4.3, using instead the solution

$$u(t, x; \gamma^*) = \begin{cases} 200xt^2(.5-x) & , \quad 0 \leq x < .3 \\ 17.143t^2(1-x) & , \quad .3 \leq x \leq 1. \end{cases}$$

We began the parameter search with the start-up guess of

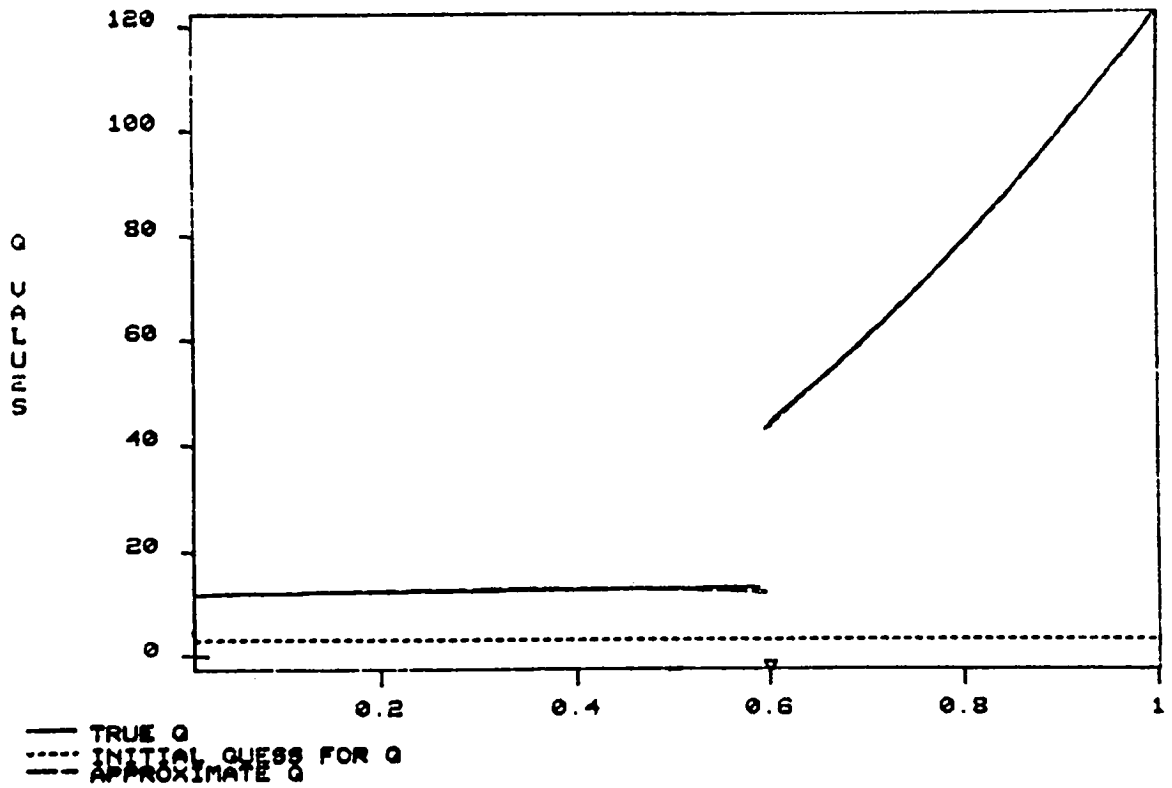


Figure 1a: Example 4.3

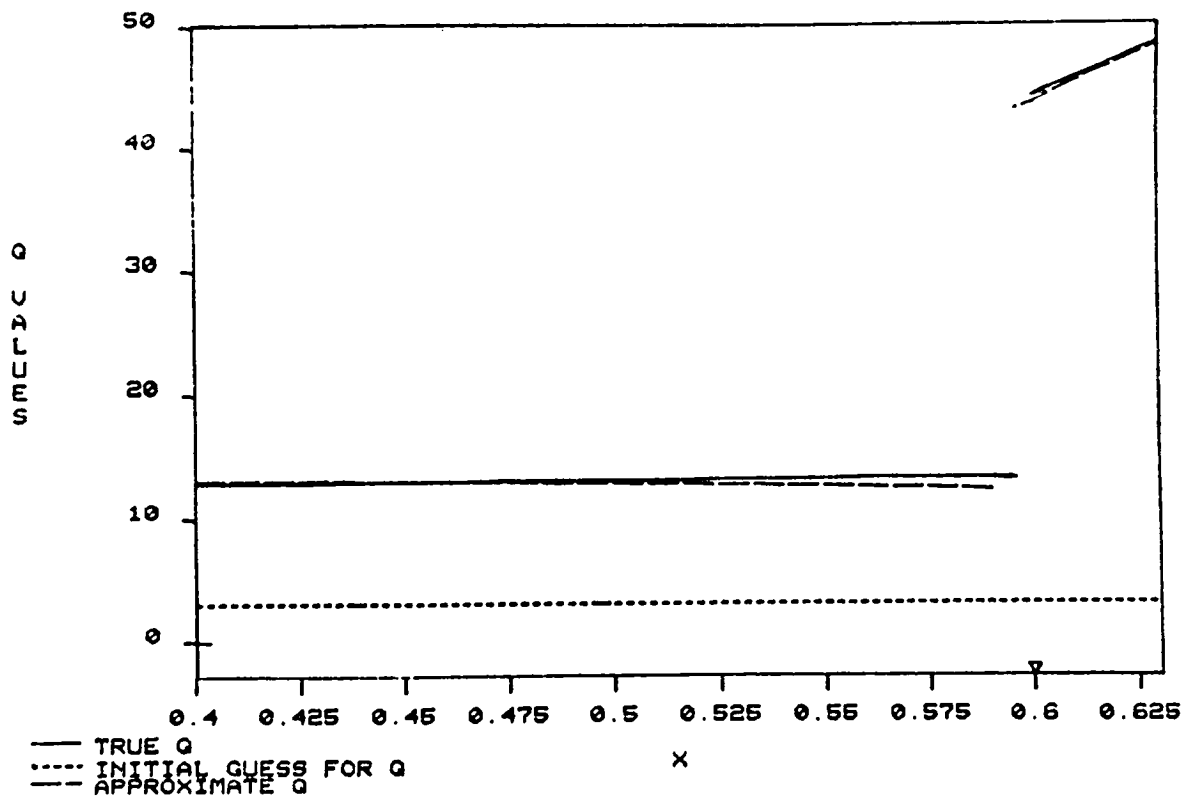


Figure 1b: Example 4.3

$$q^0 = \begin{cases} 18. & , \quad 0 \leq x < .2 \\ 48. & , \quad .2 \leq x \leq 1 \end{cases}$$

(see Figure 2a) and obtained the converged value of  $\bar{q}^{N,M}$  ( $N=24, M=1$ ) that is depicted in Figure 2b. Note that  $\bar{\xi}^{24}$  is not approaching  $\xi^*$ ; in fact, we observed that the iterates for  $\xi$  never changed from the initial guess of  $\xi^0 = .2$  (recall  $\xi^* = .3$ ) throughout the iterative process. We note that the software package did perform fairly well in its attempt to estimate the approximate functional shape of the parameter  $q$  (i.e.,  $\bar{\phi}_1^{24}, \bar{\phi}_2^{24}$  are roughly the same as  $\phi_1^*, \phi_2^*$  on the intervals  $[0., .2]$  and  $[.3, 1.]$ ; between  $x = .2$  and  $x = .3$  there is discrepancy due to the incorrect value of  $\bar{\xi}^{24}$ ). The failure of the numerical package to adequately estimate  $\xi$  may be due to some well-known limitations of the particular optimization scheme (Levenberg-Marquardt) that we chose to use with our approximation ideas: It has been our experience that difficulties sometimes arise when this minimization scheme is used to estimate more than seven or eight unknown parameters. (In this example there are 9 unknowns --  $\xi$  and 4 coefficients each in the cubic spline representations for  $\phi_1$  and  $\phi_2$ .) That the numerical package was able to estimate 9 parameters (and, in particular,  $\xi$ ) in the last example may be due to the fact that the difference between  $q(\xi^-)$  and  $q(\xi^+)$  in that example is greater, making the accurate placement of  $\xi$  more critical.

We were able to overcome the difficulties we encountered in this example by taking the following steps: First, observing that  $\xi$  had not changed at all while  $\phi_1$  and  $\phi_2$  appeared to have converged to reasonable values, we restarted the iterative process holding  $\phi_1 \equiv \bar{\phi}_1^{24}$  and  $\phi_2 \equiv \bar{\phi}_2^{24}$  fixed while iterating only on  $\xi$ . As is seen in Figure 2c, the converged value of  $\xi$  obtained using this approach is  $\bar{\xi}^{24} = .299 \approx .3 = \xi^*$ . Finally, we "re-tuned" the coefficients in the spline

Figures 2a-2d: Example 4.4

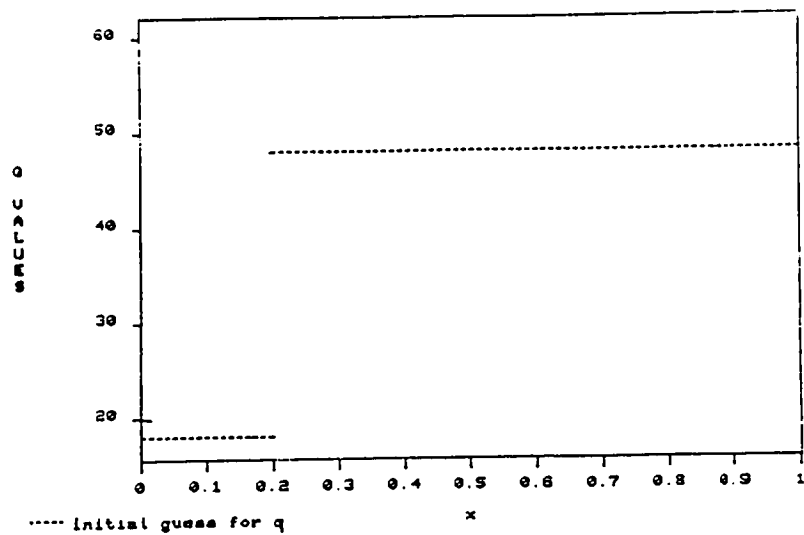


Figure 2a

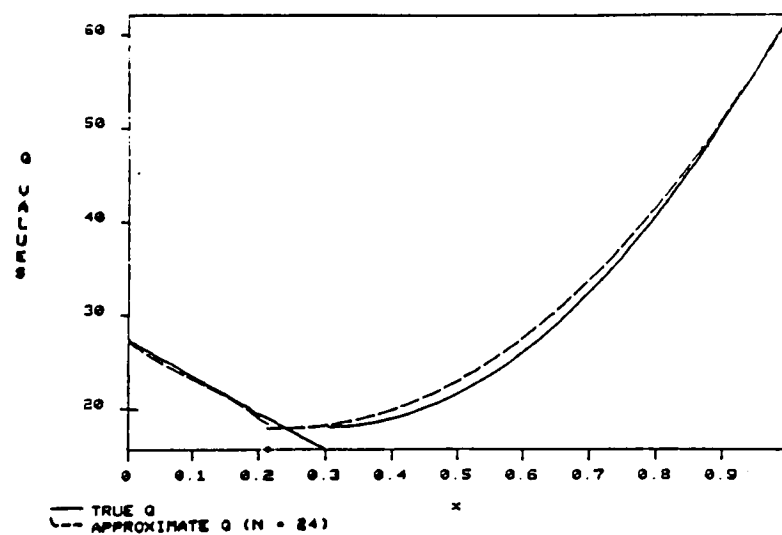


Figure 2b

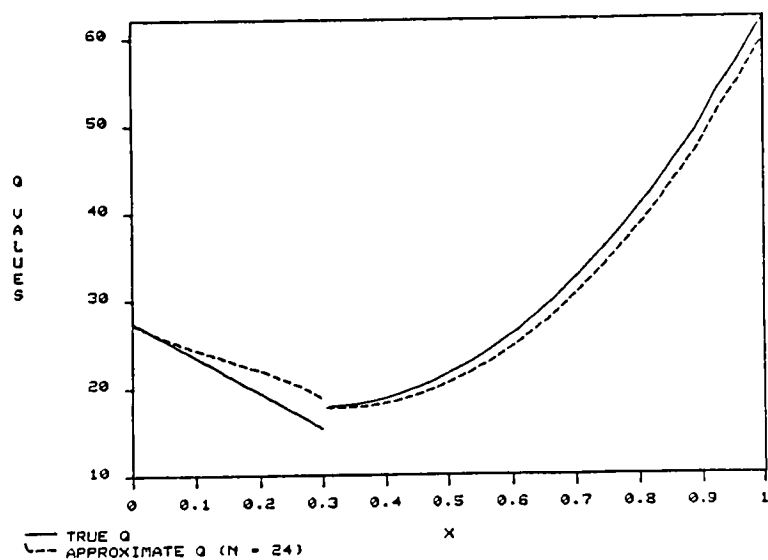


Figure 2c

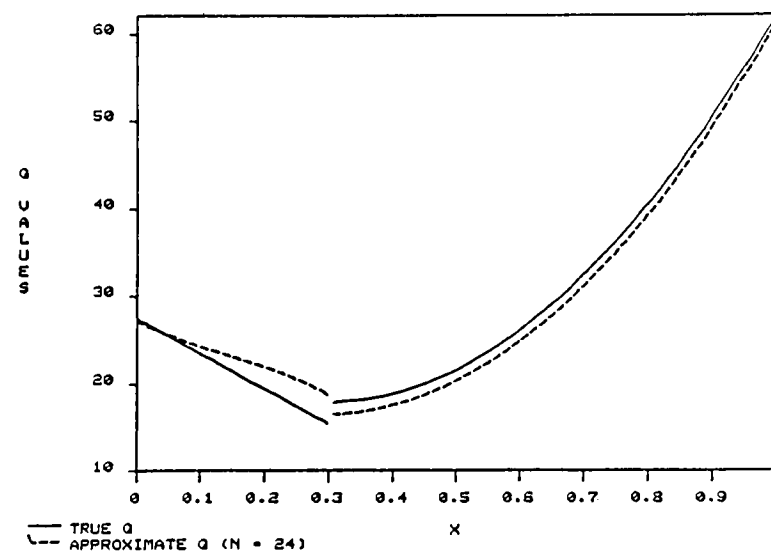


Figure 2d

expansions for  $\phi_1$  and  $\phi_2$  by iterating once again on those coefficients, this time holding  $\xi$  fixed at the new value of  $\bar{\xi}^{24}$  and using as start-up values for  $\phi_1, \phi_2$  the converged values obtained from the first iterative process (i.e.,  $\bar{\phi}_1^{24}, \bar{\phi}_2^{24}$ ). The result is shown in Figure 2d. This somewhat adaptive algorithm to estimate accurately all unknown parameters is a common approach taken (often of necessity) when real data is used in connection with model-building applications (see, for example, [11]).

Finally, we remark that a drawback of our approximation framework is that we must specify the number of discontinuities in advance of the estimation process. Fortunately, it is possible to overestimate and underestimate this number and still obtain useful information. This will be the focus of our last two examples.

Example 4.5. We repeat Example 4.2 except that we assume throughout that  $q$  is discontinuous at only one point (while two discontinuities are actually present in  $q^*$ ); we also allow spatial variation in  $\phi_1$  and  $\phi_2$  and approximate using cubic splines. An initial guess for  $q$  and a converged estimate  $\bar{q}^{N,M}$  ( $N=24, M=1$ ) are depicted in Figure 3 where it is interesting to note that the initial guess of  $\xi^0 = .4$  converges to a value close to that of the true (second) discontinuity,  $\xi_2^* = .6$ . In addition, to the right of this point the estimated shape of  $q$  begins to approximate the constant function  $\phi_3^*$ , while to the left of that point the rapidly increasing estimated shape gives an indication that we have underestimated the number of discontinuities present.

Example 4.6. We repeat Example 4.1, except that now we overestimate the number of discontinuities in  $q$ . We assume throughout that  $q = \phi_1 + H_{\xi_1} \phi_2 + H_{\xi_2} \phi_3$  where

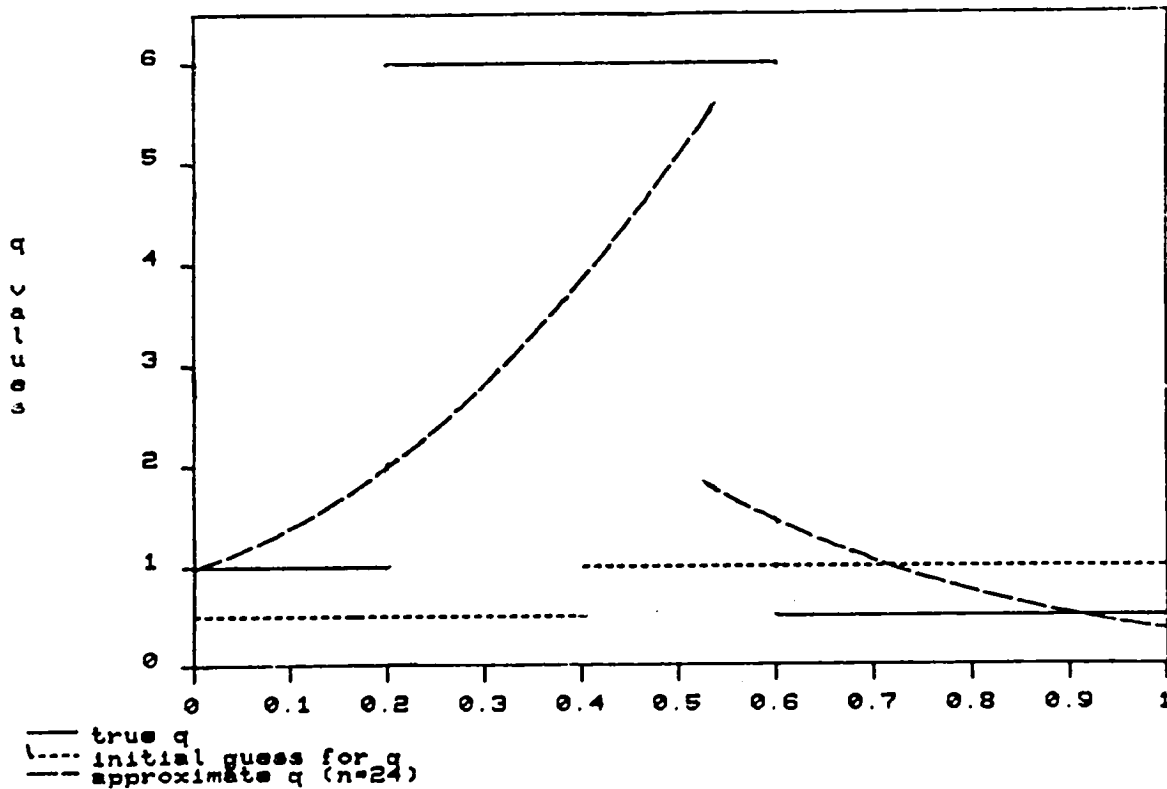


Figure 3: Example 4.5

$\phi_1$ ,  $\phi_2$ , and  $\phi_3$  are constants. For an initial guess of

$$q^0 = \begin{cases} 25. & , \quad 0 \leq x < .5 \\ 5. & , \quad .5 \leq x < .7 \\ 20. & , \quad .7 \leq x \leq 1 \end{cases} ,$$

we obtained ( $N=8$  , 291 CP seconds)

$$\bar{q}^8 = \begin{cases} 14.95 & , \quad 0. \leq x \leq .503 \\ 14.99 & , \quad .503 \leq x \leq .600 \\ 50.05 & , \quad .600 \leq x \leq 1 \end{cases} ;$$

repeating the same example but with a different initial guess,

$$q^0 = \begin{cases} .001 & , & 0 \leq x < .333 \\ .001 & , & .333 \leq x < .667 \\ .001 & , & .667 \leq x \leq 1 \end{cases} ,$$

we observed the following converged values ( $N=16$ )

$$\bar{q}^{16} = \begin{cases} 2.44 & , & 0 \leq x < .0001 \\ 15.05 & , & .0001 \leq x < .6001 \\ 49.88 & , & .6001 \leq x \leq 1. \end{cases}$$

A close inspection of either result reveals that we were, in fact, able to accurately estimate  $q^*$  (as defined in Example 4.1), even though a two-discontinuity approximation structure was incorrectly used throughout.

Remark 4.1 We note that all examples presented here involve polynomial or piecewise-polynomial state/parameter functions; from this one might conclude that such a polynomial structure is needed in order to effectively apply our spline-based methods. In fact this is not the case, as we have seen in a number of test examples (see [4], [7], [26] for a number of examples in the context of several applications).

## 5. Concluding Remarks

In the above presentation we have given a convergence theory for algorithms for the special problem of estimating discontinuous functional coefficients in parabolic systems. We are currently working to further develop and refine these ideas and to extend the theory to other applications, e.g. hyperbolic (seismic) equations and higher order (elastic beam) systems. In particular, we are studying an approximation framework that imposes the continuity condition (2.5) on approximate solutions  $u^N$  as well as on the original solution  $u$ . Our investigations in this direction involve making further (parameter dependent) modifications of the spline-based basis elements described in Section 3; we are also studying a completely different approach that involves the "tau-Legendre" ideas that have been successfully applied in [8] to (discontinuous coefficient) hyperbolic systems. We are also working to develop a related theory for two-dimensional domains, although for obvious reasons this is not simply a trivial extension of the ideas presented thus far.

Finally, we note that we have not addressed here the problem of "identifiability" of parameters, or the ill-posed nature of parameter estimation or "inverse" problems in general. These important and difficult questions arise in all parameter estimation problems and are not special difficulties associated with the particular problem under consideration here. The reader is referred to Remarks 3.3 and 3.4 of [14] for further comments regarding the nature of this ill-posedness, of nonuniqueness and the importance of initial guesses for parameters, and the general unavailability of convergence rates for approximate parameters  $\{q^{N,M}\}$ .

## 6. Acknowledgment

The author would like to express appreciation to Prof. H. T. Banks for numerous insightful discussions during the course of this work.



## REFERENCES

- [1] R. A. Adams, Sobolev Spaces, Academic Press, N. Y., 1975.
- [2] H. T. Banks, J. A. Burns, and E. M. Cliff, Parameter estimation and identification for systems with delays, SIAM J. Control & Optimization 19 (1981), 791-828.
- [3] H. T. Banks and J. M. Crowley, Parameter estimation for distributed systems arising in elasticity, Proc. Symposium on Engineering Sciences and Mechanics, (National Cheng Kung University, Tainan, Taiwan, Dec. 28-31, 1981), 158-177; LCDS Tech. Rep. 81-24, Brown University, November, 1981.
- [4] H. T. Banks and J. M. Crowley, Parameter identification in continuum models, LCDS Rep. M-83-1, Brown University, March, 1983; Proc. Amer. Cont. Conf., San Francisco, June, 1983, 997-1001.
- [5] H. T. Banks, J. M. Crowley and K. Kunisch, Cubic spline approximation techniques for parameter estimation in distributed systems, IEEE Trans. Auto. Control 28 (1983), 773-786.
- [6] H. T. Banks, P. L. Daniel (Lamm) and E. S. Armstrong, A spline-based parameter and state estimation technique for static models of elastic surfaces, ICASE Rep. No. 82-25, NASA LRC, Hampton, Va., June, 1983.
- [7] H. T. Banks, P. L. Daniel (Lamm), and E. S. Armstrong, Spline-based estimation techniques for parameters in elliptic distributed systems, Proc. Fourth VPI & SU/AIAA Symposium on Dynamics and Control of Large Structures (Blacksburg, June 1983), to appear; LCDS Rep. No. 83-22, Brown University, June, 1983.
- [8] H. T. Banks, K. Ito, and K. A. Murphy, Computational methods for estimation of parameters in hyperbolic systems, Proc. Conf. on Inverse Scattering: Theory and Application (Tulsa, May 16-18, 1983), SIAM, Philadelphia, 1983.
- [9] H. T. Banks and P. Kareiva, Parameter estimation techniques for transport equations with application to population dispersal and tissue bulk flow models, J. Math. Biology 17 (1983), 253-273.
- [10] H. T. Banks, P. M. Kareiva, P. K. Daniel Lamm, Estimation techniques for transport equations, Proc. Int'l. Conf. on Mathematics in Biology and Medicine (Bari, July 18-22, 1983), to appear; LCDS Tech. Rep. No. 83-23, Brown University, July, 1983.
- [11] H. T. Banks, P. M. Kareiva, and P. K. Daniel Lamm, Estimation of temporally and spatially varying coefficients in models for insect dispersal, LCDS Tech. Rep. No. 83-14, Brown University, June, 1983; J. Math. Biology, submitted.

- [12] H. T. Banks and K. Kunisch, An approximation theory for nonlinear partial differential equations with applications to identification and control, SIAM J. Control and Optimization 20 (1982), 815-849.
- [13] H. T. Banks and P. K. Daniel Lamm, Estimation of delays and other parameters in nonlinear functional differential equations, SIAM J. Control and Optimization 21 (1983), 895-915.
- [14] H. T. Banks and Patricia Daniel Lamm, Estimation of variable coefficients in parabolic distributed systems, LCDS Rep. No. 82-22, Sept. 1982, Brown University, Providence, RI 02912; IEEE Trans. Auto. Control, to appear.
- [15] J. M. Crowley, Numerical Methods of Parameter Identification for Problems Arising in Elasticity, Ph.D. Dissertation, Brown University, May, 1982.
- [16] P. L. Daniel (Lamm), Spline-based approximation methods for the identification and control of nonlinear functional differential equations, Ph.D. Dissertation, Brown University, Providence, RI, June, 1981.
- [17] R. E. Ewing, Determination of Coefficients in Reservoir Simulation, Numerical Treatment of Inverse Problems in Differential and Integral Equations, P. Deufhard and E. Hairer, editors, Birkhäuser, Boston, 1983, 206-226.
- [18] R. E. Ewing, The mathematics of reservoir simulation, Ch. I, SIAM Frontiers in Appl. Math. 1 (1984), to appear.
- [19] D. Gilbarg and N. S. Trudinger, Elliptic Partial Differential Equations of Second Order, Springer-Verlag, New York, 1977.
- [20] J. A. Goldstein, Semigroups of Operators and Abstract Cauchy Problems, Lecture Notes, Tulane University, 1970.
- [21] D. Gottlieb and S. Orszag, Numerical Analysis of Spectral Methods: Theory and Applications, Vol. 26, CBMS-NSF Reg. Conf. Series in Applied Math., SIAM, Philadelphia, 1977.
- [22] K. Kunisch, Identification and estimates of parameters in abstract Cauchy problems, Preprint. No. 11 (1981), Technical University of Graz.
- [23] K. Kunisch and L. W. White, The parameter estimation problem for parabolic equations in multidimensional domains in the presence of point evaluations, Preprint No. 17 (1983), Technical University of Graz.
- [24] K. Kunisch and L. W. White, Parameter estimation for elliptic equations in multidimensional domains with point and flux observations, Preprint No. 19 (1983), Technical University of Graz.
- [25] K. A. Murphy, A Spline-based Approximation Method for Inverse Problems for a Hyperbolic System Including Unknown Boundary Parameters, Ph.D. Dissertation, Brown Univ., Providence, RI, May, 1983.

- [26] K. A. Murphy, A spline-based approximation method for inverse problems for a hyperbolic system including unknown boundary parameters, Proc. 22nd IEEE Conf. on Decision and Control, San Antonio, Dec., 1983, to appear.
- [27] A. Pazy, Semigroups of Linear Operators and Applications to Partial Differential Equations, Applied Mathematical Sciences 44, Springer-Verlag, New York, 1983.
- [28] W. Rudin, Functional Analysis, McGraw-Hill, New York, 1973.
- [29] M. H. Schultz, Spline Analysis, Prentice-Hall, Englewood Cliffs, N. J., 1973.
- [30] M. F. Wheeler,  $L_\infty$  estimates of optimal orders for Galerkin methods for one-dimensional second order parabolic and hyperbolic equations, SIAM J. Numerical Analysis 10 (1973), 908-913.

1. Report No. NASA CR-172301		2. Government Accession No.		3. Recipient's Catalog No.	
4. Title and Subtitle  Estimation of Discontinuous Coefficients in Parabolic Systems: Applications to Reservoir Simulation				5. Report Date February 1984	
				6. Performing Organization Code	
7. Author(s) Patricia Daniel Lamm				8. Performing Organization Report No. 84-7	
9. Performing Organization Name and Address Institute for Computer Applications in Science and Engineering Mail Stop 132C, NASA Langley Research Center Hampton, VA 23665				10. Work Unit No.	
				11. Contract or Grant No. NASI-17130, NASI-16394	
12. Sponsoring Agency Name and Address National Aeronautics and Space Administration Washington, D.C. 20546				13. Type of Report and Period Covered Contractor report	
				14. Sponsoring Agency Code	
15. Supplementary Notes Additional Support: NSF Grant MCS-8200883, NASA Grant NAG-1-258.  Langley Technical Monitor: Robert H. Tolson Final Report					
16. Abstract  We present spline-based techniques for estimating spatially varying parameters that appear in parabolic distributed systems (typical of those found in reservoir simulation problems). In particular, we discuss the problem of determining discontinuous coefficients, estimating both the functional shape and points of discontinuity for such parameters. In addition, our ideas may also be applied to problems with unknown initial conditions and unknown parameters appearing in terms representing external forces. Convergence results and a summary of numerical performance of the resulting algorithms are given.					
17. Key Words (Suggested by Author(s)) parameter estimation discontinuous coefficients parabolic systems			18. Distribution Statement 64 Numerical Analysis  Unclassified-Unlimited		
19. Security Classif. (of this report) Unclassified	20. Security Classif. (of this page) Unclassified	21. No. of Pages 55	22. Price A04		



