

# A WEIGHTED LEAST SQUARES METHOD FOR FIRST-ORDER HYPERBOLIC SYSTEMS

D. G. ZEITOUN\*, J. P. LAIBLE AND G. F. PINDER

*College of Engineering and Mathematics, 101 Votey Building, University of Vermont, Burlington, VT 05403, U.S.A.*

## SUMMARY

The paper presents a generalization of the classical  $L^2$ -norm weighted least squares method for the numerical solution of a first-order hyperbolic system. This alternative least squares method consists of the minimization of the weighted sum of the  $L^2$  residuals for each equation of the system. The order of accuracy of global conservation of each equation of the system is shown to be inversely proportional to the weight associated with the equation. The optimal relative weights between the equations are then determined in order to satisfy global conservation of the energy of the physical system.

As an application of the algorithm, the shallow water equations on an irregular domain are first discretized in time and then solved using Laplace modification and the proposed least squares method.

KEY WORDS Weighted least squares Hyperbolic system of equations Shallow water equations Newton–Raphson method

## 1. INTRODUCTION

While there are a number of different least squares formulations for solving boundary value problems, they may be generally classified into three groups.

In the first group the partial differential equation is transformed into an optimal control problem via the introduction of a state function, the solution of a given state equation. This equation is generally discretized using a Galerkin finite element approximation. The least squares solution algorithm employs a conjugate gradient method. Conservation of mass and satisfaction of boundary conditions are assured by requiring the state vector to belong to a suitable Sobolev space.<sup>1,2</sup>

In the second group the least squares approximation is combined with the Galerkin finite element method to obtain convergent finite element approximations. An example of this type of method is the Galerkin least squares method developed by Hughes *et al.*<sup>3</sup> In this approach the least squares form of the residual is added to the Galerkin method in order to stabilize the Galerkin method without degrading accuracy.

In the third group the mathematical form of the approximate solution (trial solution) is first chosen and then the norm of the corresponding residual is minimized by least squares.<sup>4</sup> This is the approach addressed in the present paper.

A collocation  $L^2$  least squares formulation is proposed herein which takes into account in an optimal way the relative weights identified with the different equations of the system. While the methodology developed in this paper may be applied to various types of partial differential equations, herein the solution of a first-order hyperbolic system generally denoted as the shallow water equations is considered.

---

\*Present address: TAHAL Consulting Eng. Ltd., P.O. Box 111700-6111, Tel-Aviv, Israel.

After a brief review in Section 2 of numerical methods for solving first-order systems of partial differential equations, basic definitions and mathematical properties of the weighted least squares method are discussed in Sections 3 and 4. The absence of global conservation properties of the classical least squares formulation leads to a new method proposed in Section 5 and applied to the solution of the shallow water equations in Section 6.

## 2. NUMERICAL METHODS FOR FIRST-ORDER SYSTEMS

Let  $\Omega$  be a bounded domain in the  $(x, y)$  plane with boundary  $\Gamma$ .

The vector  $\mathbf{n} = (n_x, n_y)$  is the outward unit normal to  $\Gamma$ .

The first-order system of coupled equations considered in this work may be written as

$$\mathbf{L}\mathbf{u} \equiv \mathbf{A} \frac{\partial \mathbf{u}}{\partial x} + \mathbf{D} \frac{\partial \mathbf{u}}{\partial y} + \mathbf{C}\mathbf{u} = \mathbf{f}, \quad (1)$$

with the boundary condition

$$\mathbf{B}\mathbf{u} \equiv (\mathbf{A}n_x + \mathbf{D}n_y - \mathbf{M})\mathbf{u} = \mathbf{g}. \quad (2)$$

$\mathbf{A}$ ,  $\mathbf{D}$ ,  $\mathbf{C}$  and  $\mathbf{M}$  are  $p \times p$  matrix-valued functions, the column vector forcing term is  $\mathbf{f}^T = (f_1, f_2, \dots, f_p)$ ,  $\mathbf{g}^T = (g_1, g_2, \dots, g_{pb})$  is the column vector of applied boundary conditions and the unknown column vector is  $\mathbf{u}^T = (u_1, u_2, \dots, u_p)$  ( $T$  denotes transposition).  $\mathbf{L}$  represents the linear differential operator defined in the domain  $\Omega$  and  $\mathbf{B}$  represents the boundary differential operator.

Sufficient conditions for problem (1), (2) to have a unique strong solution have been obtained for symmetric positive systems in the sense of Friedrichs.<sup>5</sup>

Numerical approximations of pure hyperbolic problems have been extensively studied. Different schemes based on finite difference and Galerkin finite element approximations have been proposed for the solution of linear and non-linear hyperbolic equations and first-order systems (see Reference 6 for a review).

A perceived disadvantage of the least squares method compared with the Galerkin method is the higher-order finite element continuity requirement. However, for basis functions defined on rectangular subspaces, such as in the case presented herein, these continuity conditions are easily accommodated and higher-order accuracy and function continuity are achieved. These higher-order continuity constraints may be circumvented by expressing higher-order equations as sets of lower-order equations, although additional field variables are thereby introduced.

The standard Galerkin method leads to 'central-type' discrete operators which exhibit oscillatory behaviour on practical meshes. A way to avoid this difficulty is to select a test function for the convective term that explicitly accommodates the directional property of the hyperbolic propagation (upwinding). Such methods are often denoted as Petrov-Galerkin methods. The main drawback of this approach is the absence of a generally applicable systematic procedure for the selection of the test function. Recently the Lax-Wendroff scheme has been used for developing the Taylor-Galerkin method. In addition, the method of characteristics has been used for developing the Galerkin method.

Finite element methods for solving first-order systems which are symmetric and positive in the sense of Friedrichs have been proposed by Lesaint and Raviart.<sup>7</sup> For symmetric systems arising out of the heat equation, Aziz and Liu<sup>8</sup> proposed a weighted least squares solution method. In their work a unique weight associated with the boundary conditions was determined in order to reduce the error estimates.

For mixed methods based on splines, the least squares theory for second-order elliptic systems developed in Reference 9 specifies the optimal weights in terms of the spline mesh size. Other first-order systems have also been solved numerically using a weighted least squares formulation where

optimal weights are found theoretically. For example, Wendland<sup>10</sup> formulates a least squares method for strong elliptic systems satisfying the Lopatinsky conditions for the boundary operator. For elliptic operators in the sense of Douglis and Nirenberg, wherein the operator satisfies the supplementary condition for the operator  $\mathbf{L}$  and the complementary conditions for the boundary operator  $\mathbf{B}$ , a weighted least squares method has been proposed by Aziz *et al.*<sup>11</sup>

A difficulty appearing in the implementation of the method, and one that is addressed in this paper, is the determination of the weights associated with the least squares method for a general first-order system. In most numerical implementations these weights must be specified by the analyst.<sup>4</sup> Without some rational criteria for selecting these weights, the efficiency of the method may be seriously affected even if other computational objectives are achieved.

The formulation proposed herein is similar to the least squares finite element method proposed by Jiang and Carey for linear and non-linear hyperbolic systems.<sup>12-14</sup> However, the method differs in as much as collocation points and weighting functions are introduced.

As noted earlier, the least squares method is validated on the two-dimensional shallow water equations. These equations are often used to obtain flow fields necessary for pollution problems and transport in shallow estuaries. A special advantage of using the collocation least squares method for solving the shallow water equations is the ability to solve the equations in an irregular domain with a completely orthogonal computational mesh.<sup>15,16</sup> This significantly enhances both the accuracy and computation time.

When the non-linear terms appearing in these equations are linearized, they reduce to a system of first-order equations. Classical collocation least squares has good noise control characteristics compared with the Galerkin method when solving the shallow water equations. However, one of the difficulties inhibiting its widespread use is that the accuracy of the solution depends upon the weights appearing in the formulation.<sup>16</sup>

### 3. A WEIGHTED LEAST SQUARES METHOD

A least squares formulation is proposed herein which takes into account the relative weights between the different first-order equations appearing in the system. The least squares formulation is presented in terms of the matrix of differential operators  $\mathbf{L}$  and  $\mathbf{B}$  defined in equations (1) and (2). These equations may be written as

$$\mathbf{L}\mathbf{u} = \mathbf{f} \quad \text{in } \Omega, \quad \mathbf{B}\mathbf{u} = \mathbf{g} \quad \text{on } \Gamma. \quad (3)$$

In order to formulate the method proposed, let us now define the functional spaces and their associated norms. In the following we define the working space  $V_g$  as the  $p$ -product of the Sobolev space  $H^1(\Omega)$ :

$$V_g = \underbrace{H^1(\Omega) \times H^1(\Omega) \times \cdots \times H^1(\Omega)}_{p \text{ times}} = [H^1(\Omega)]^p, \quad (4)$$

$$H^1(\Omega) = \left\{ v \in L^2(\Omega); \frac{\partial v}{\partial x} \in L^2(\Omega); \frac{\partial v}{\partial y} \in L^2(\Omega) \right\}.$$

Notice that each element  $\mathbf{u} \in V_g$  is a vector of  $p$  components  $u_i$ ,  $i = 1, 2, \dots, p$ .

If for any  $\mathbf{u} \in V_g$  and  $\mathbf{v} \in V_g$ ,  $\mathbf{u} \cdot \mathbf{v}$  denotes the classical scalar product on  $\mathbb{R}^p$ , an  $L^2$  scalar product may be defined on the domain  $\Omega$  and the boundary  $\Gamma$  as

$$(\mathbf{u}, \mathbf{v}) = \int_{\Omega} \mathbf{u} \cdot \mathbf{v} d\Omega, \quad [\mathbf{u}, \mathbf{v}] = \int_{\Gamma} \mathbf{u} \cdot \mathbf{v} d\Gamma. \quad (5)$$

In the following,  $\|\cdot\|_{\Omega}$  denotes the norm associated with  $(\cdot, \cdot)$  and  $\|\cdot\|_{\Gamma}$  denotes the norm associated with  $[\cdot, \cdot]$ .

The  $L^2$  weighted least squares formulation proposed in this work may be separated into two groups.

*Group 1.  $L^2$  continuous least squares*

This formulation corresponds to a minimization problem with the choice of an  $L^2(\Omega)$ -norm for the interior domain and an  $L^2(\Gamma)$ -type norm for the boundary conditions:

$$\begin{aligned} \text{Min}_{(\mathbf{v} \in V_b)} T(\mathbf{v}) &= (\mathbf{L}\mathbf{v} - \mathbf{f}, \mathbf{W}(\mathbf{L}\mathbf{v} - \mathbf{f})) + [\mathbf{B}\mathbf{v} - \mathbf{g}, \mathbf{P}(\mathbf{B}\mathbf{v} - \mathbf{g})] \iff \\ &= \text{Min}_{(\mathbf{v} \in V_b)} \sum_{i=1}^p w_i \int_{\Omega} \left( \sum_{j=1}^p L_{ij} v_j - f_i \right)^2 d\Omega + \sum_{i=1}^p p_i \int_{\Gamma} \left( \sum_{j=1}^{pb} B_{ij} v_j - g_i \right)^2 d\Gamma, \quad (6) \end{aligned}$$

where  $\mathbf{W}$  and  $\mathbf{P}$  represent the diagonal matrices of positive weights associated respectively with each equation of the system ( $w_i$ ) and each boundary condition ( $p_i$ ).

In practice, problem (6) is never solved on a functional space of infinite dimension. A finite-dimensional approximation of problem (6) is generally solved by restricting the above functional to be on a finite-dimensional subspace of  $V_g$ . One such approach is based on finite element approximations.

With this approach the discretization procedure may be formulated as follows.

The closed domain  $\Omega$  is approximated by a polygonal domain  $\Omega_h$  with a standard triangulation  $\mathcal{T}_h$  of  $\Omega_h$ , i.e.  $\mathcal{T}_h$  is a set of finite elements  $E$ .<sup>17</sup>

The Sobolev space  $H^1(\Omega)$  is then approximated on  $\mathcal{T}_h$  by

$$H_h^1(\Omega_h) = \{v_h | v_h \in C^0(\bar{\Omega}_h); v_{h|E} \in P_1, \forall E \in \mathcal{T}_h\}, \quad (7)$$

where  $P_1$  is the space of polynomials in two variables of degree less than or equal to one. In the treatment of parabolic problems written in terms of second-order operators, the least squares procedure requires a space of polynomials which are at least  $P_2$ . In this case the Hermite interpolation functions are generally used for the interpolation over a single element.<sup>6,17</sup>

The space  $V_h$ , the approximation of  $V_g$ , is then defined as the  $p$ -product of  $H_h^1(\Omega_h)$ . Each component  $v_{hj}$  of the vector function  $\mathbf{v}_h \in V_h$  may be written in terms of the global interpolation function defined as

$$\Phi^j(\mathbf{x}) = \begin{pmatrix} \phi_1^j(\mathbf{x}) \\ \vdots \\ \phi_{N_p^j}^j(\mathbf{x}) \end{pmatrix}, \quad (8)$$

$$\phi_i^j(\mathbf{x}) \in C^0(\bar{\Omega}_h), \quad (9)$$

where  $N_p^j$  is the total number of unknowns for the function  $v_j(\mathbf{x})$ .

Formally, any function  $v_{hj} \in H_h^1(\Omega_h)$  may be written in terms of the vector of global interpolation functions and the vector of unknown parameters,  $\mathbf{a}^j$ , as

$$v_{hj}(\mathbf{x}) = \sum_{E \in \mathcal{T}_h} \sum_{i=1}^{N_p^j} \phi_{iE}^j(\mathbf{x}) a_i^j = [\Phi^j(\mathbf{x})] \mathbf{a}^j, \quad (10)$$

where  $N_T$  is the number of unknowns on the element  $E$ . The independent space variable of the domain is represented by  $\mathbf{x}$ . The  $L^2$  approximate continuous least squares formulation consists of the minimization of a functional  $I_c(\mathbf{a}, \mathbf{W}, \mathbf{P})$  containing a weighted sum of the interior and boundary residuals with respect to the unknown vectors  $\mathbf{a}^i$ .

The interior residual  $R_{Li}$  corresponding to the  $i$ th equation of the system and the boundary residuals  $R_{Bi}$  corresponding to the  $i$ th boundary condition are defined by

$$\begin{aligned} R_{Li}(\mathbf{a}, \mathbf{x}) &= \sum_{j=1}^p L_{ij} v_{hj}(\mathbf{x}) - f_i, \quad \mathbf{x} \text{ in } \Omega, \\ R_{Bi}(\mathbf{a}, \mathbf{x}) &= \sum_{j=1}^{pb} B_{ij} v_{hj}(\mathbf{x}) - g_i, \quad \mathbf{x} \text{ on } \Gamma. \end{aligned} \quad (11)$$

Consider now the global unknown vector  $\mathbf{a}$  formed by the set of all the vectors  $\mathbf{a}^j$ ,  $j = 1, \dots, p$ . The objective function  $I_c(\mathbf{a}, \mathbf{W}, \mathbf{P})$  to minimize may be defined as

$$I_c(\mathbf{a}, \mathbf{W}, \mathbf{P}) = \sum_{i=1}^p w_i \int_{\Omega} [R_{Li}(\mathbf{a}, \mathbf{x})]^2 d\Omega + \sum_{i=1}^{pb} p_i \int_{\Gamma} [R_{Bi}(\mathbf{a}, \mathbf{x})]^2 d\Gamma. \quad (12)$$

The approximate  $L^2$  continuous least squares formulation corresponding to equation (6) may be defined as

$$\text{Min}_{(\mathbf{v}_h \in V_h)} I_c(\mathbf{a}, \mathbf{W}, \mathbf{P}). \quad (13)$$

In this general framework an error analysis of the least squares approximation has been conducted for particular problems in fluid dynamics.<sup>1,4,5,18</sup> The choice of a suitable norm  $\|\cdot\|_{\Omega}$  on  $V_g$  depends on the type of differential operator. For example, the choice of an  $H^1$ -norm instead of the classical  $L^2$ -norm stabilized the least squares solution of non-linear hyperbolic problems which generate shocks.<sup>13</sup>

#### Group 2. $L^2$ collocation (or discrete) least squares

The discrete formulation corresponding to the continuous least squares method, also called collocation least squares, consists of selecting a series of collocation points inside the domain and on the boundary and minimizing the function

$$I_d(\mathbf{a}, \mathbf{W}, \mathbf{P}) = \sum_{i=1}^p w_i \sum_{l=1}^k c_l [R_{Li}(\mathbf{a}, \mathbf{x}_l)]^2 d\Omega + \sum_{i=1}^{pb} p_i \sum_{l=k+1}^m c_l [R_{Bi}(\mathbf{a}, \mathbf{x}_l)]^2 d\Gamma. \quad (14)$$

The points  $\mathbf{x}_l$  for  $l = 1, \dots, k$  correspond to the interior points and for  $l = k+1, \dots, m$  to the boundary points of  $\Phi$ .<sup>4</sup> The weights  $c_l$  are associated with the collocation points. One of the difficulties appearing in the implementation of the method is the determination of the weights  $w_i$ ,  $i = 1, \dots, p$ , the weights  $p_i$ ,  $i = 1, \dots, p$ , and the weights  $c_l$ ,  $l = 1, \dots, m$ .

The addition of the integral of the boundary residual may result in a solution vector with large interior residuals. This scaling problem may degrade the satisfaction of the governing balance laws. This is particularly true when the weights appearing in the matrix  $\mathbf{P}$  are very large.

Basically, the three types of weights introduced ( $w_i$ ,  $p_i$  and  $c_i$ ) are each of a different nature. The first and second types of weights,  $w_i$ ,  $i = 1, \dots, p$ , and  $p_i$ ,  $i = 1, \dots, p_b$ , determine the relative contributions of the interior and boundary residuals. When a system of equations is solved using an  $L^2$  least squares method, these weights will balance the different contributions of each equation and each boundary condition. These weights are generally specified by the analyst.

The third type of weights,  $c_i$ ,  $i = 1, \dots, m$ , represent the distribution of weights among the collocation points. These weights may be determined by using optimal criteria for numerical integration, such as Gaussian quadrature. For a collocation least squares method using finite element discretization, the use of Gaussian points and their weights leads to an optimal truncation error. More sophisticated criteria have been proposed by the authors.<sup>19</sup>

In the present work we concentrate on the first and second types of weights and a general methodology is developed for computation of the optimal values of these weights.

#### 4. PROPERTIES OF THE $L^2$ LEAST SQUARES

In the least squares literature two different methods have been proposed for computing the numerical solution of the minimization problem (problem (13) here). The first consists of using an unconstrained optimization algorithm for the quadratic functional  $I_c$  or  $I_d$ . Classical methods such as conjugate gradient or quasi-Newton generally require the computation of both a functional and its Jacobian. The second approach consists of solving the normal equation corresponding to a first-order necessary condition for a minimum of the functional.

In the present work we deduce the normal equation from the multidimensional minimization problem (equation (6)) rather than from the finite-dimensional problem as is usually done. Then we discuss the basic properties of the least squares method.

The following theorem gives a necessary condition for a minimum for the least squares and approximate least squares methods.

##### *Theorem 1*

If  $\tilde{\mathbf{u}}$ , resp.  $\mathbf{u}_h$ , is the vector solution of the minimization problem (6), resp. (13), then for all  $\mathbf{v} \in V_g$ , resp.  $\mathbf{v}_h \in V_h$ .

$$(\mathbf{L}\tilde{\mathbf{u}} - \mathbf{f}, \mathbf{W}(\mathbf{L}\mathbf{v})) + [\mathbf{B}\tilde{\mathbf{u}} - \mathbf{g}, \mathbf{P}(\mathbf{B}\mathbf{v})] = 0, \quad (15)$$

$$(\mathbf{L}\mathbf{u}_h - \mathbf{f}, \mathbf{W}(\mathbf{L}\mathbf{u}_h)) + [\mathbf{B}\mathbf{u}_h - \mathbf{g}, \mathbf{P}(\mathbf{B}\mathbf{v}_h)] = 0. \quad (16)$$

*Proof.* Consider the least squares functional  $J(\mathbf{v}) = (\mathbf{L}\mathbf{v} - \mathbf{f}, \mathbf{W}(\mathbf{L}\mathbf{v} - \mathbf{f})) + [\mathbf{B}\mathbf{v} - \mathbf{g}, \mathbf{P}(\mathbf{B}\mathbf{v} - \mathbf{g})]$ . If  $\tilde{\mathbf{u}}$  is a solution of equation (15), then  $\forall \mathbf{v} \in V_g$

$$\lim_{t \rightarrow 0} \frac{J(\tilde{\mathbf{u}} + t\mathbf{v}) - J(\tilde{\mathbf{u}})}{t} \geq 0. \quad (17)$$

The Taylor expansion of  $J(\tilde{\mathbf{u}} + t\mathbf{v})$  may be written as

$$J(\tilde{\mathbf{u}} + t\mathbf{v}) = J(\tilde{\mathbf{u}}) + t\{(\mathbf{L}\tilde{\mathbf{u}} - \mathbf{f}, \mathbf{W}(\mathbf{L}\mathbf{v})) + [\mathbf{B}\tilde{\mathbf{u}} - \mathbf{g}, \mathbf{P}(\mathbf{B}\mathbf{v})]\} + O(t^2).$$

Then for all  $\mathbf{v} \in V_g$  and for  $t > 0$  the above limit leads to

$$(\mathbf{L}\tilde{\mathbf{u}} - \mathbf{f}, \mathbf{W}(\mathbf{L}\mathbf{v})) + [\mathbf{B}\tilde{\mathbf{u}} - \mathbf{g}, \mathbf{P}(\mathbf{B}\mathbf{v})] \geq 0. \quad (18)$$

Taking now  $-\mathbf{v}$  instead of  $\mathbf{v}$  in the last inequality, we see that we obtain (15).

The same proof may be developed for the approximate least squares problem (equation (16)).  $\square$

Equation (16) may be written as a linear system of algebraic equations with unknown vector  $\mathbf{a}$ . The solution vector  $\mathbf{u}_h$  and the weighting function vector  $\mathbf{v}_h$  may be expressed in terms of their components as

$$u_{hj}(\mathbf{x}) = [\Phi^j(\mathbf{x})]\mathbf{a}^j, \quad v_{hj}(\mathbf{x}) = [\Phi^j(\mathbf{x})]\mathbf{d}^j. \quad (19)$$

Equation (16) may be written in terms of the components of the differential operators  $\mathbf{L}$  and  $\mathbf{B}$  as

$$\begin{aligned} \sum_{i=1}^p w_i \int_{\Omega} \left( \sum_{j=1}^p L_{ij}([\Phi^j(\mathbf{x})])\mathbf{a}^j - f_i \right) \left( \sum_{l=1}^p L_{il}([\Phi^l(\mathbf{x})])\mathbf{d}^l \right) d\Omega \\ + \sum_{i=1}^{pb} p_i \int_{\Gamma} \left( \sum_{j=1}^{pb} B_{ij}([\Phi^j(\mathbf{x})])\mathbf{a}^j - g_i \right) \left( \sum_{l=1}^{pb} B_{il}([\Phi^l(\mathbf{x})])\mathbf{d}^l \right) d\Gamma = 0. \end{aligned} \quad (20)$$

Let us assume that for any  $j = 1, \dots, p$  the dimension of the column vector  $\mathbf{a}^j$  is independent of  $j$  and equals  $D$ . Introduce now the  $N^D \times pN^D$  matrices  $\mathbf{V}$  and the two multipliers  $\mathbf{L}_i(\mathbf{x})$  and  $\mathbf{B}_i(\mathbf{x})$  defined as

$$\mathbf{a}^j = \mathbf{V}\mathbf{a}, \quad \mathbf{L}_i(\mathbf{x}) = \sum_{j=1}^p L_{ij}([\Phi^j(\mathbf{x})])\mathbf{V}, \quad \mathbf{b}_i(\mathbf{x}) = \sum_{j=1}^{pb} B_{ij}([\Phi^j(\mathbf{x})])\mathbf{V}. \quad (21)$$

With this notation equation (16) may be written in matrix form as

$$\left[ \left( \sum_{i=1}^p w_i \mathbf{A}_i \right) \mathbf{a} - \sum_{i=1}^p w_i \mathbf{f}_i \right] \cdot \mathbf{d} + \left[ \left( \sum_{i=1}^{pb} p_i \mathbf{B}_i \right) \mathbf{a} - \sum_{i=1}^{pb} p_i \mathbf{g}_i \right] \cdot \mathbf{d} = 0, \quad (22)$$

where the matrices  $\mathbf{A}_i$  and  $\mathbf{B}_i$  and the vectors  $\mathbf{f}_i$  and  $\mathbf{g}_i$  are defined for continuous least squares as

$$\mathbf{A}_i = \int_{\Omega} \mathbf{L}_i^L(\mathbf{x})\mathbf{L}_i(\mathbf{x})d\mathbf{x}, \quad \mathbf{B}_i = \int_{\Gamma} \mathbf{b}_i^T(\mathbf{x})\mathbf{b}_i(\mathbf{x})d\mathbf{x}, \quad \mathbf{f}_i = \int_{\Omega} \mathbf{L}_i^T(\mathbf{x})\mathbf{f}_i d\mathbf{x}, \quad \mathbf{g}_i = \int_{\Gamma} \mathbf{b}_i^T(\mathbf{x})\mathbf{g}_i d\mathbf{x} \quad (23)$$

and for collocation least squares as

$$\mathbf{A}_i = \sum_{j=1}^k v_j \mathbf{L}_i^T(\mathbf{x}_j)\mathbf{L}_i(\mathbf{x}_j), \quad \mathbf{B}_i = \sum_{j=k+1}^m \mathbf{b}_i^T(\mathbf{x}_j)\mathbf{b}_i(\mathbf{x}_j), \quad \mathbf{f}_i = \sum_{j=1}^k \mathbf{L}_i^T(\mathbf{x}_j)\mathbf{f}_i, \quad \mathbf{g}_i = \sum_{j=k+1}^m \mathbf{b}_i^T(\mathbf{x}_j)\mathbf{g}_i. \quad (24)$$

Equation (22) is valid for any vector  $\mathbf{d}$  and may therefore be written in the classical form of the normal equation, i.e.

$$\left( \sum_{i=1}^p w_i \mathbf{A}_i + \sum_{i=1}^{pb} p_i \mathbf{B}_i \right) \mathbf{a} = \sum_{i=1}^p w_i \mathbf{f}_i + \sum_{i=1}^{pb} p_i \mathbf{g}_i. \quad (25)$$

Without loss of generality, equation (25) may be scaled such that the first weights  $w_1$  will be equal to one. Defining the new weights  $\epsilon_i = 1/w_i$ ,  $i = 2, \dots, p$ , and  $v_j = 1/w_j$ ,  $j = 1, \dots, p_b$ , equation (25) may now be expressed as

$$\left( \mathbf{A}_1 + \sum_{i=2}^p \frac{1}{\epsilon_i} \mathbf{A}_i + \sum_{i=1}^{p_b} \frac{1}{v_i} \mathbf{B}_i \right) \mathbf{a} = \mathbf{f}_1 + \sum_{i=1}^p \frac{1}{\epsilon_i} \mathbf{f}_i + \sum_{i=1}^{p_b} \frac{1}{v_i} \mathbf{g}_i. \quad (26)$$

The choice of the optimal relative weightings depends on the criterion of optimality defined. For the numerical least squares method presented in this contribution, these criteria may be separated into four categories:

- (a) accuracy of the numerical scheme
- (b) numerical stability of the linear system resulting from the normal equation (25)
- (c) satisfaction of the boundary conditions
- (d) conservation of the mass balance.

The dependence of the weighting on the overall accuracy of the numerical scheme is still a subject of research. Error estimates using different norms on suitable Sobolev spaces have been derived by several authors for the penalty method (see e.g. Reference 9). Compared with our formulation, the classical penalty method corresponds to a system with two equations and a single weight. The optimal weight is chosen for the reduction of the error estimation. This weight depends on the type of partial differential equation, the discretization parameter and the type of boundary conditions. Such analysis shows that a better theoretical error estimate is obtained when the weight is reduced.

It is well known that for a large weight the condition number of the normal equation resulting from the penalty least squares method depends linearly on this weight.<sup>1</sup> According to this criterion, one has to increase the weight for a good stability of the normal equation.

Thus the weighting strategy for criterion (a) is the opposite of that for criterion (b).

Among the basic properties we may require from a numerical method, global conservation of the system is of primary importance.<sup>3</sup> This property depends on the physical problem to be solved and on the type of differential operator. It may be stated that the flux, the energy or the forces will be conserved globally inside  $\Omega$ . For the first-order hyperbolic system considered here, the global conservation property may be written for each equation of the system as follows:

$$\text{for } i = 1, \dots, p, \quad \int_{\Omega} \left( \sum_{j=1}^p L_{ij} v_j - f_i \right) d\Omega = 0. \quad (27)$$

The global conservation property is not always satisfied by the numerical solution of the weighted least squares formulation presented herein. On the basis of Lagrangian multiplier functions, the order of accuracy of the global conservation property is analysed. The following theorem presents an error estimate for the global conservation property for each equation of the first-order system.

### Theorem 2

If  $\mathbf{a}$  is the vector of dimension  $pN^D$  which is the solution of equation (26), then we have the following.

- (i) For the continuous least squares method the following properties hold:

$$\begin{aligned} \text{for } i = 1, \dots, p, \quad & \int_{\Omega} \left( \sum_{j=1}^p L_{ij} u_{hj} - f_i \right) d\Omega = H_i(\mathbf{x}) \epsilon_i + O(\epsilon_i^2), \\ \text{for } i = 1, \dots, pb, \quad & \sum_{j=1}^{pb} B_{ij} u_{hj} - g_i = D_i v_i + O(v_i^2). \end{aligned} \quad (28)$$



(ii) For the collocation least squares method the following properties hold:

$$\begin{aligned} \text{for } i = 1, \dots, p, \quad \sum_{l=1}^k v_l \left( \sum_{j=1}^p L_{ij} u_{hj}(\mathbf{x}_l) - f_i \right) d\Omega &= H_i(\mathbf{x}) \epsilon_i + O(\epsilon_i^2), \\ \text{for } i = 1, \dots, pb \quad \text{for } l = k+1, \dots, m, \quad \sum_{j=1}^{pb} B_{ij} u_{hj}(\mathbf{x}_l) - g_i &= D_i(\mathbf{x}) v_i + O(v_i^2). \end{aligned} \quad (29)$$

Here the function  $H_i(\mathbf{x})$  does not depend on  $\epsilon_i$  and  $D_i$  does not depend on  $v_i$ .

*Proof.* The proof is developed in the Appendix. □

For first-order systems a rigorous study of the dependence of the weights on the accuracy of our least squares formulation is outside the scope of this contribution. However, Theorem 2 points out the difficulty in choosing the relative weights appearing in the least squares functional between the different equations.

In terms of global mass balance conservation, the weighting strategy is to increase all the weight appearing in equations (28) and (29). However, in terms of conditioning of the matrix, the weighting strategy is to reduce these weights. For a given equation (*i*) the choice of a small value for one of the  $\epsilon_i$ ,  $i = 2, \dots, p$ , will improve the global conservation property for this equation. A good strategy should be to choose large weights for each equation, but unfortunately this will cause ill conditioning in the linear system of equations (26).

So far the goal of Theorem 2 and the above discussion was to justify the need for the definition of a strategy of optimal weighting.

In the following a criterion of optimality based on energy balance considerations is introduced for the computation of these weights.

The optimal weights are those for which the corresponding least squares solution respects more accurately the mechanical energy balance of the system. This balance equation corresponds to the natural proportions between the different equations.

For the first-order hyperbolic system considered here, one may require a conservation property of the type

$$\sum_{i=1}^p \int_{\Omega} \left( \sum_{j=1}^p L_{ij} v_j - f_i \right) v_i = 0, \quad (30)$$

which may be written in a functional form as

$$\Phi_a(\mathbf{v}) = 0. \quad (31)$$

The property of global conservation is implicitly verified in the Galerkin method but is not respected by the formulation (6).

The basic idea of the new least squares formulation proposed is to determine the optimal set of weights  $w_i$ ,  $i = 1, \dots, p$ , and  $p_i$ ,  $i = 1, \dots, p$ , in order that the least squares solution will respect approximately the global conservation law  $\Psi_a(\mathbf{v}) = 0$ . In the case where  $\Psi_a(\mathbf{v})$  is a positive function, the determination of the optimal set of weights may be achieved by minimizing the function  $\Psi_a(\mathbf{v})$  with respect to the weights.

## 5. OPTIMIZATION ALGORITHM

## 5.1. General methodology

The basic steps of the methodology proposed are as follows.

Consider after  $k$  iterations the set of known weights  $\mathbf{W}^k$  and  $\mathbf{P}^k$ .

*Step 1.* Determine  $\mathbf{u}_h(\mathbf{x}, \mathbf{W}^k, \mathbf{P}^k) = \text{ArgMin } J(\mathbf{v})$ , the solution of the minimization problem (6).

*Step 2.* For a positive function  $\Psi_a(\mathbf{v})$ , determine the new matrices with positive coefficients,  $\mathbf{W}^{(k+1)}$  and  $\mathbf{P}^{(k+1)}$ , which are obtained via solution of the minimization problem with respect to the weights:

$$\text{Min} \Psi_a(\mathbf{u}_h, \mathbf{W}^{(k+1)}, \mathbf{P}^{(k+1)}) \quad (w_i > 0; p_i > 0; i = 1, \dots, p).$$

*Step 3*

If  $|\Psi_a(\mathbf{u}_h, \mathbf{W}^{(k+1)}, \mathbf{P}^{(k+1)})| < \epsilon$  then END

If not

Set  $\mathbf{W}^k = \mathbf{W}^{(k+1)}$  and  $\mathbf{P}^k = \mathbf{P}^{(k+1)}$

GO TO Step 1.

In the next section the algorithm developed for the solution of the shallow water equations is presented.

## 5.2. Methods and procedures

Here we describe the process for determining the weights used in the present least squares method. An associate norm derived from an auxiliary finite element formulation is first computed. Then the weights are determined in order to force the auxiliary norm to a minimum. In the following we will need to distinguish between the discretized form of the differential equations as derived from the least squares method and the discretized form of the differential equations as derived by some other finite element procedure. We will denote the discretized least squares formulation equation (25) as

$$\mathbf{A}\mathbf{x} - \mathbf{f} = \mathbf{0}, \quad (32)$$

where

$$\mathbf{A} = \sum_{i=1}^p (w_i \mathbf{A}_i + p_i \mathbf{B}_i), \quad \mathbf{x} = \mathbf{a}, \quad \mathbf{f} = \sum_{i=1}^p (w_i \mathbf{f}_i + p_i \mathbf{g}_i), \quad (33)$$

and the discretized auxiliary formulation as

$$\mathbf{A}_a \mathbf{x}_a - \mathbf{f}_a = \mathbf{0}. \quad (34)$$

One possible auxiliary norm may be defined as

$$\Phi_a = [\mathbf{A}_a \mathbf{x}(\mathbf{S}) - \mathbf{f}_a]^T [\mathbf{A}_a \mathbf{x}(\mathbf{S}) - \mathbf{f}_a], \quad (35)$$

where  $\mathbf{x}(\mathbf{S})$  represents the least squares solution for a given set of weights  $\mathbf{S} = (s_1, s_2, \dots, s_{2n}) = (w_1, w_2, \dots, w_n, p_1, \dots, p_n)$  and T denotes the transposed vector. Notice that  $\mathbf{x}(\mathbf{S})$  is not the solution of the auxiliary formulation. Since  $\mathbf{A}$  and  $\mathbf{f}$  are dependent on the weights used in the least squares formulation,  $\mathbf{x}$  is also a function of the weights. Our minimization problem may thus be

stated as

$$\text{Min}_{s_1, s_2, \dots, s_{2n}} \Phi_a(\mathbf{x}(\mathbf{S})) \text{ subject to } \mathbf{A}(\mathbf{S})\mathbf{x}(\mathbf{S}) - \mathbf{f}(\mathbf{S}) = \mathbf{0}, \quad s_i > 0. \quad (36)$$

Here we detail a procedure to solve this problem and contrast this approach with the well-known Newton–Raphson and quasi-Newton methods. Consider a first-order Taylor expansion of  $\mathbf{x}(\mathbf{S})$  around  $\mathbf{S}$ :

$$\Phi_a = [\mathbf{A}_a(\mathbf{x} + \nabla_s \mathbf{x} \cdot (\Delta \mathbf{S})) - \mathbf{f}_a] \cdot [\mathbf{A}_a(\mathbf{x} + \nabla_s \mathbf{x} \cdot (\Delta \mathbf{S})) - \mathbf{f}_a] + O(\|\Delta \mathbf{S}\|^2), \quad (37)$$

where  $\nabla_s \mathbf{x}$  represents the gradient vector of the vector  $\mathbf{x}$  with respect to the weights (also called the sensitivity matrix). Applying the minimization criterion  $\partial \Psi_a / \partial \Delta \mathbf{S} = \mathbf{0}$  as in a least squares problem, we obtain

$$[\mathbf{G}^T \mathbf{A}_a^T \mathbf{A}_a \mathbf{G}](\Delta \mathbf{S}) = -[\mathbf{G}^T \mathbf{A}_a^T] \{\mathbf{A}_a \mathbf{x} - \mathbf{f}_a\}, \quad (38)$$

where  $\mathbf{G} = \nabla_s \mathbf{x}$  is the sensitivity matrix. This matrix is obtained by differentiation with respect to each weight  $s_i$  of the least squares equation:

$$\frac{\partial}{\partial s_i} [\mathbf{A} \mathbf{x} - \mathbf{f}]. \quad (39)$$

Expanding this last equation and solving for  $\partial \mathbf{x} / \partial s_i$  yields

$$\mathbf{g}_i = \frac{\partial \mathbf{x}}{\partial s_i} = \mathbf{A}^{-1} \left\{ \frac{\partial \mathbf{f}}{\partial s_i} - \left[ \frac{\partial \mathbf{A}}{\partial s_i} \right] \mathbf{x} \right\}. \quad (40)$$

The sensitivity matrix  $\mathbf{G}$  is thus defined as  $\mathbf{G} = [\mathbf{g}_1 | \mathbf{g}_2 | \mathbf{g}_3 | \dots | \mathbf{g}_n]$ , where  $n$  is the number of weights. The steps of the procedure are as follows.

1. Assume a set of weights  $\mathbf{S}$ .
2. Solve  $\mathbf{A} \mathbf{x} - \mathbf{f} = \mathbf{0}$ .
3. Employ equation (40) to obtain  $\mathbf{G}$ .
4. Form and solve equation (38) to obtain  $\Delta \mathbf{S}$ .
5. If  $\|\Delta \mathbf{S}\| < \epsilon_{\text{tol}}$ , end; else  $\mathbf{S}^{\text{new}} = \mathbf{S}^{\text{old}} + \alpha \{\Delta \mathbf{S}\}$  and go to step 2.

Alternatively one may expand  $\Psi_a$  around the vector  $\mathbf{S}_0$  by Taylor's theorem to obtain

$$\Phi_a = \Phi_a(\mathbf{S}_0) + \nabla_s(\Phi_a) \cdot \{\Delta \mathbf{S}\} + \frac{1}{2} \{\Delta \mathbf{S}\} \cdot \mathbf{H}_s \{\Delta \mathbf{S}\} + O(\|\Delta \mathbf{S}\|^2), \quad (41)$$

where  $\mathbf{H}_s$  is the Hessian matrix of  $\Psi_a$  with respect to the weights and  $\Delta \mathbf{S} = \mathbf{S} - \mathbf{S}_0$ . The minimum of  $\Psi_a$  is obtained from  $\partial \Psi_a / \partial \Delta \mathbf{S} = \mathbf{0}$ . This yields

$$\mathbf{H}_s \{\Delta \mathbf{S}\} = -\nabla_s(\Phi_a) = \mathbf{G}^T \mathbf{A}_a^T [\mathbf{A}_a \mathbf{x} - \mathbf{f}_a] \quad (42)$$

This last equation is the basis for the second-order Newton–Raphson methods. It can be shown that the exact expression for the second-order derivatives is

$$\mathbf{H}_w = [\mathbf{G}^T \mathbf{A}_a^T \mathbf{A}_a \mathbf{G}] + \left[ \left[ \frac{\partial \mathbf{G}}{\partial \mathbf{S}} \right] \right] [\mathbf{A}_a^T (\mathbf{A}_a \mathbf{x} - \mathbf{f}_a)], \quad (43)$$

where  $[[\partial \mathbf{G} / \partial \mathbf{S}]]$  is a third-order tensor. For the problem under consideration this tensor can be determined analytically but not without considerable computational expense. To improve on the inverse of the Hessian matrix, one can employ the BFGS methods readily available in the IMSL Fortran Library.

## 6. HYPERBOLIC SHALLOW WATER EQUATIONS

Here we discuss a specific case of the general methodology defined in Section 4. The specific equations used are the shallow water equations, which may be written in scalar form as

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} + v \frac{\partial u}{\partial y} + \frac{C_b}{H} u - fv + g \frac{\partial \zeta}{\partial x} - \frac{\tau_x}{H} = 0, \quad (44)$$

$$\frac{\partial v}{\partial t} + u \frac{\partial v}{\partial x} + v \frac{\partial v}{\partial y} + \frac{C_b}{H} v + fu + g \frac{\partial \zeta}{\partial y} - \frac{\tau_y}{H} = 0, \quad (45)$$

$$\frac{\partial \zeta}{\partial t} + \frac{\partial(Hu)}{\partial x} + \frac{\partial(Hv)}{\partial y} - r = 0, \quad (46)$$

$$\zeta - \zeta_p = 0, \quad (47)$$

$$l_x u + l_y v = v_n, \quad (48)$$

$$H = h + \zeta. \quad (49)$$

The first two equations are the  $x$ - and  $y$ -momentum equations respectively. The third equation is the fluid continuity equation. These equations are valid over the region  $\Omega$ . The fourth equation is a prescribed surface elevation to be enforced on a portion  $\Phi_1$  of the boundary and the fifth equation is a prescribed normal flow to be applied on a portion  $\Phi_2$  of the boundary. The last equation defines the total depth. For problems with moving boundaries this equation becomes a third-type boundary condition. Here we will restrict our attention to problems of constant geometry. In these equations  $x$  and  $y$  are Cartesian co-ordinates,  $t$  is time,  $g$  is gravity,  $f$  is the Coriolis parameter,  $C_b$  is the bottom friction parameter,  $u$  and  $v$  are  $x$  and  $y$  vertically integrated velocities respectively,  $\zeta$  is the surface elevation,  $h$  is the water depth from mean sea level,  $\zeta_p$  is the prescribed surface elevation,  $v_n$  is the prescribed normal velocity,  $l_x$  and  $l_y$  are the direction cosines of the outwardly directed unit normal on the boundary,  $\tau_x$  and  $\tau_y$  represent the  $x$ - and  $y$ -components of the wind shear stress respectively and  $r$  is the fluid source. The bottom friction parameter may take on various forms and is generally dependent on  $u$  and  $v$ .

These equations in their full non-linear form have been solved by the least squares method on irregular domains developed by Laible and Pinder.<sup>15</sup> Since we wish to use a test problem that has an analytic solution, we will focus on the linearized form of these equations.

Starting with the scalar equations, we now seek to bring the shallow water equations into the form of equations (1) and (2). First we note that the continuity equations can be expanded in terms of  $h$  and  $\xi$ . Inserting the equality  $H = h + \xi$  into the continuity equation, the following continuity equation may be obtained:

$$\frac{\partial \xi}{\partial t} + u \frac{\partial \xi}{\partial x} + v \frac{\partial \xi}{\partial y} + \left( \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} \right) \xi + h \frac{\partial u}{\partial x} + h \frac{\partial v}{\partial y} + \frac{\partial h}{\partial x} u + \frac{\partial h}{\partial y} v - \gamma = 0. \quad (50)$$

This equation is now in terms of the first-order derivatives of the unknowns  $u$ ,  $v$  and  $\xi$ .

Equations (44)–(46) are non-linear. We now seek to obtain a linearized form. Here we introduce a modification (similar to the idea of Laplace modification used to solve problems with time-dependent coefficients). This modification is applied to all the non-linear terms of equations (44)–(46). Suppose we select some time-invariant characteristic values for  $u$ ,  $v$  and  $h$  denoted as  $U$ ,  $V$  and  $H_i$ . These values must generally be greater than any anticipated values of  $u$ ,  $v$  and  $h$  respectively. The modification is accomplished by writing an identity equation for each of equations (44)–(46) containing the non-linear

terms. To illustrate the process, we will consider equation (44) and deduce the remaining ‘Laplace-modified’ equations. The identity for equation (44) may be written as

$$U \frac{\partial u}{\partial x} + V \frac{\partial u}{\partial y} + \frac{C_b}{Hi} u = U \frac{\partial u}{\partial x} + V \frac{\partial u}{\partial y} + \frac{C_b}{Hi} u. \quad (51)$$

The wind stress term  $\tau_x/H$  and  $C_b$  are taken here to be independent of  $u$ ,  $v$  and  $\xi$ , although they could also be similarly treated.

If we now subtract this identity from equation (44), one obtains

$$\frac{\partial u}{\partial t} + U \frac{\partial u}{\partial x} + V \frac{\partial u}{\partial y} + \frac{C_b}{H} u - fv + g \frac{\partial \zeta}{\partial x} \quad (52)$$

$$= (U - u) \frac{\partial u}{\partial x} + (V - v) \frac{\partial u}{\partial y} + \left( \frac{C_b}{Hi} - \frac{C_b}{H} \right) u + \frac{\tau_x}{H}. \quad (53)$$

This process itself does not introduce any numerical approximations. In the numerical formulation, however, the left-hand side now contains a linear operator with constant coefficients. In the solution process values of  $u$ ,  $v$  and  $H$  are assumed known (from a previous time step or by extrapolation from previous values). Therefore the right-hand side will contain known values. After solution for  $u$ ,  $v$  and  $\zeta$  the right-hand side is updated and the system solved again. Iteration within a time step is carried out until some norm of the variation of  $u$ ,  $v$  and  $\zeta$  meets a convergence tolerance. In the following we will denote the assumed or extrapolated values for  $u$ ,  $v$  and  $\zeta$  on the right-hand side as  $u^*$ ,  $v^*$  and  $\zeta^*$  (also  $H^* = h + \zeta^*$ ).

Applying this modification to equations (45) and (46), we finally have

$$\frac{\partial u}{\partial t} + U \frac{\partial u}{\partial x} + V \frac{\partial u}{\partial y} + \frac{C_b}{H} u - fv + g \frac{\partial \zeta}{\partial x} = \hat{F}_u, \quad (54)$$

$$\frac{\partial v}{\partial t} + U \frac{\partial v}{\partial x} + v \frac{\partial v}{\partial y} + \frac{C_b}{H} v + fu + g \frac{\partial \zeta}{\partial y} = \hat{F}_v, \quad (55)$$

$$\frac{\partial \zeta}{\partial t} + U \frac{\partial \zeta}{\partial x} + v \frac{\partial \zeta}{\partial y} + \left( \frac{\partial U}{\partial x} + \frac{\partial V}{\partial y} \right) \zeta + h \frac{\partial u}{\partial x} + h \frac{\partial v}{\partial y} + \frac{\partial h}{\partial x} u + \frac{\partial h}{\partial y} v = \hat{F}_\zeta, \quad (56)$$

where

$$\hat{F}_u = (U - u^*) \frac{\partial u^*}{\partial x} + (V - v^*) \frac{\partial u^*}{\partial y} + \left( \frac{C_b}{H} - \frac{C_b^*}{H^*} \right) u^* + \frac{\tau_x}{H^*} \quad (57)$$

$$\hat{F}_v = (U - u^*) \frac{\partial v^*}{\partial x} + (V - v^*) \frac{\partial v^*}{\partial y} + \left( \frac{C_b}{H} - \frac{C_b^*}{H^*} \right) v^* + \frac{\tau_x}{H^*} \quad (58)$$

$$\hat{F}_\zeta = (U - u^*) \frac{\partial \zeta^*}{\partial x} + (V - v^*) \frac{\partial \zeta^*}{\partial y} + \left[ \left( \frac{\partial U}{\partial x} + \frac{\partial V}{\partial y} \right) - \left( \frac{\partial u^*}{\partial x} + \frac{\partial v^*}{\partial y} \right) \right] \zeta^* + \gamma. \quad (59)$$

We now introduce the time discretization scheme. The local time derivative term is approximated by a backward difference, e.g.  $\partial u / \partial t \approx (u^+ - u^-) / \Delta t$ . The variables  $u$ ,  $v$ ,  $\zeta$ ,  $u^*$ ,  $v^*$  and  $\zeta^*$  are evaluated at an intermediate time between  $t$  and  $t + \Delta t$  by the trapezoidal rule, e.g.  $u = \alpha u^+ + (1 - \alpha) u^-$ . Here  $u^+$  is

at time  $t + \Delta t$  and  $u^-$  is at time  $t$ . With these approximations we may finally collect all terms that multiply  $u$ ,  $\partial u/\partial x$  and  $\partial u/\partial y$  as in the standard form of equation (1). The vector  $\mathbf{u}$  is now taken as

$$\mathbf{u} = \begin{Bmatrix} u \\ v \\ \xi \end{Bmatrix}. \quad (60)$$

After some rearrangement we find

$$\mathbf{A}_\alpha \frac{\partial \mathbf{u}^+}{\partial x} + \mathbf{D}_\alpha \frac{\partial \mathbf{u}^+}{\partial y} + \mathbf{C}_\alpha \mathbf{u}^+ = \mathbf{f}, \quad (61)$$

where

$$\begin{aligned} \mathbf{A}_\alpha &= \begin{bmatrix} \alpha U & 0 & \alpha g \\ 0 & \alpha U & 0 \\ \alpha h & 0 & \alpha U \end{bmatrix}, & \mathbf{D}_\alpha &= \begin{bmatrix} \alpha V & 0 & 0 \\ 0 & \alpha V & \alpha g \\ 0 & \alpha h & \alpha V \end{bmatrix}, \\ \mathbf{C}_\alpha &= \begin{bmatrix} \frac{1}{\Delta t} + \frac{C_b}{H} \alpha & -\alpha f & 0 \\ \alpha f & \frac{1}{\Delta t} + \frac{C_b}{H} \alpha & 0 \\ \alpha \frac{\partial h}{\partial x} & \alpha \frac{\partial h}{\partial y} & \frac{1}{\Delta t} + \alpha \left( \frac{\partial U}{\partial x} + \frac{\partial V}{\partial y} \right) \end{bmatrix}, \\ \mathbf{F} &= \begin{Bmatrix} \hat{F}_u \\ \hat{F}_v \\ \hat{F}_\xi \end{Bmatrix} + \mathbf{A}_\beta \frac{\partial \mathbf{u}^-}{\partial x} + \mathbf{D}_\beta \frac{\partial \mathbf{u}^-}{\partial y} + \mathbf{C}_\beta \mathbf{u}^-. \end{aligned} \quad (62)$$

The matrices  $[\mathbf{A}]_\beta$ ,  $[\mathbf{D}]_\beta$  and  $[\mathbf{C}]_\beta$  are identical to  $[\mathbf{A}]_\alpha$ ,  $[\mathbf{D}]_\alpha$  and  $[\mathbf{C}]_\alpha$  except that  $\alpha$  is replaced by  $\beta = (\alpha - 1)$ .

The boundary conditions are also written in the standard form (equation (2)) as

$$\begin{bmatrix} 0 & 0 & 1 \\ n_x & n_y & 0 \end{bmatrix} \begin{Bmatrix} u^+ \\ v^+ \\ \xi^+ \end{Bmatrix} = \begin{Bmatrix} \xi \\ v_p \end{Bmatrix}.$$

### Residuals

It is now possible to express the residuals of the differential equations and boundary conditions in terms of the standard matrices  $[\mathbf{A}]$ ,  $[\mathbf{D}]$ ,  $[\mathbf{C}]$  and  $[\mathbf{M}]$ . Here we use the matrix notation dropping the subscript  $\alpha$ . To develop the residuals due to time and space discretization, we now introduce the spatial approximations. We introduce basis function expansions of our unknowns with local support over rectangular finite elements:

$$\{\mathbf{u}\} \approx [\Phi] \{\bar{\mathbf{u}}\}, \quad (63)$$

where

$$[\Phi] = \begin{bmatrix} [\phi_i] & 0 & 0 \\ 0 & [\phi_i] & 0 \\ 0 & 0 & [\phi_i] \end{bmatrix}.$$

The vector  $[\phi_i]$  contains the four basic functions for the bilinear rectangular element. The matrix  $[\epsilon]$  is thus a  $3 \times 12$  matrix and  $\{\bar{\mathbf{u}}\}$  is a  $12 \times 1$  vector of nodal values on an elemental region.

Substitution of (63) into equation (61) yields the residuals over the domain  $\Omega$ ,

$$\epsilon_{\Omega} = \begin{Bmatrix} \epsilon_u \\ \epsilon_v \\ \epsilon_{\xi} \end{Bmatrix} = \left[ [\mathbf{A}]_{\alpha} \frac{\partial \Phi}{\partial x} + [\mathbf{D}]_{\alpha} \frac{\partial \Phi}{\partial y} + [\mathbf{C}]_{\alpha} \Phi \right] \{\bar{\mathbf{u}}\} - \mathbf{f},$$

and over the boundary  $\Gamma$ ,

$$\epsilon_{\Gamma} = \begin{Bmatrix} \epsilon_{b\xi} \\ \epsilon_{bu} \end{Bmatrix} = [\mathbf{M}][\Phi]\{\bar{\mathbf{u}}\}^+ - \{\mathbf{g}\}.$$

Collectively we may write the residuals as

$$\epsilon_{\Omega} = \begin{Bmatrix} \epsilon_u \\ \epsilon_v \\ \epsilon_{\xi} \\ \epsilon_{b\xi} \\ \epsilon_{bu} \end{Bmatrix} = \begin{bmatrix} [A]_{\alpha} \frac{\partial \Phi}{\partial x} + [D]_{\alpha} \frac{\partial \Phi}{\partial y} + [C]_{\alpha} \Phi \\ [M][\Phi] \end{bmatrix} \{\bar{\mathbf{u}}\}^+ - \begin{Bmatrix} f \\ g \end{Bmatrix}$$

or equivalently

$$\{\epsilon\} = [\mathbf{LB}^+]\{\bar{\mathbf{u}}\}^+ - [\mathbf{LB}^-]\{\bar{\mathbf{u}}\}^- - \begin{Bmatrix} \hat{f} \\ \hat{g} \end{Bmatrix}, \quad (64)$$

where  $\{\bar{\mathbf{u}}\}^+$  and  $\{\bar{\mathbf{u}}\}^-$  are the vectors of nodal values of the variables at two different time levels and  $[\mathbf{LB}^+]$  and  $[\mathbf{LB}^-]$  are the numerical counterparts of the differential operator matrix  $[\mathbf{LB}]$ .

The total squared weighted residual is given by

$$E^2 = \sum_{\Omega, \Gamma} \{\epsilon\}^T [\mathbf{S}] \{\epsilon\}. \quad (65)$$

Substitution of equation (64) into equation (65) and application of the minimization criterion  $\partial E^2 / \partial \{\mathbf{u}\}^+ = \mathbf{0}$  leads to

$$[\mathbf{A}]^+ \{\mathbf{u}\}^+ - \{\mathbf{F}\} = \mathbf{0}, \quad (66)$$

where

$$[\mathbf{A}]^+ = \sum_{\Omega, \Gamma} [\mathbf{LB}^+]^T [\mathbf{S}] [\mathbf{LB}^+], \quad (67)$$

$$[\mathbf{A}]^- = \sum_{\Omega, \Gamma} [\mathbf{LB}^+]^T [\mathbf{S}] [\mathbf{LB}^-], \quad (68)$$

$$\{\mathbf{P}\} = \sum_{\Omega, \Gamma} [\mathbf{LB}^+] [\mathbf{S}] \begin{Bmatrix} \hat{f} \\ \hat{g} \end{Bmatrix}, \quad \{\mathbf{F}\} = [\mathbf{A}]\{\mathbf{u}\}^- + \{\mathbf{P}\}. \quad (69)$$

### 6.1. Test problem description

The test problem is shown in Figure 1. The quarter-ring domain has a quadratically varying bathymetry defined by  $h = H_0 r^2$ ,  $H_0 = h_1 / r_1^2$ ,  $h_1 = 16$  m. The surface elevation is prescribed at the outer boundary ( $r = r^2$ ) and all other boundaries have a zero-normal-flow condition. The steady state numerical and analytical solutions were obtained for an  $x$ -directed constant wind loading. The velocity

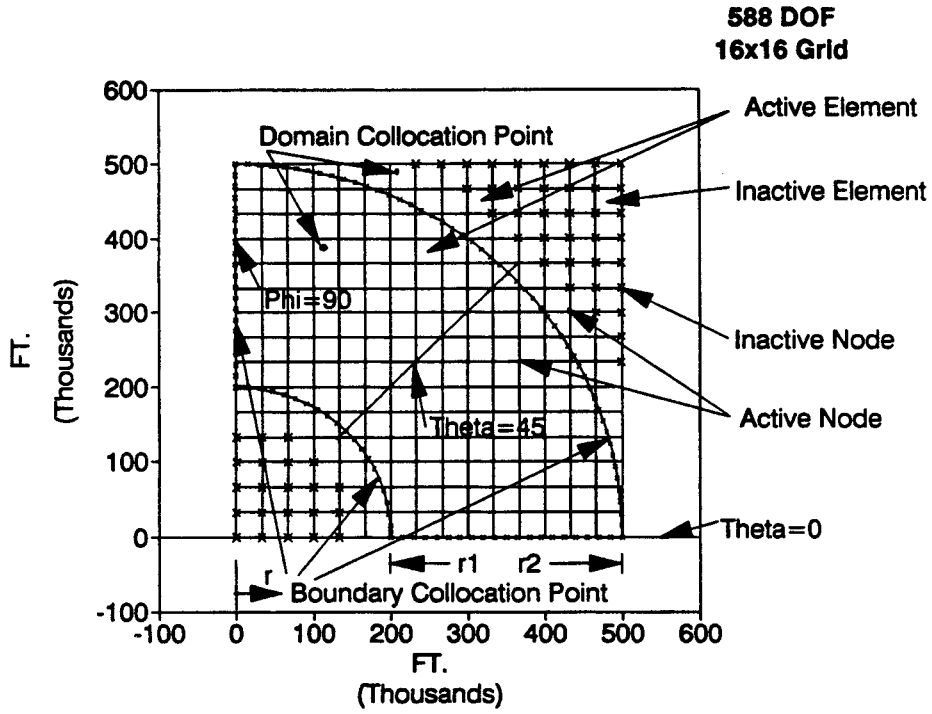


Figure 1. Test grid

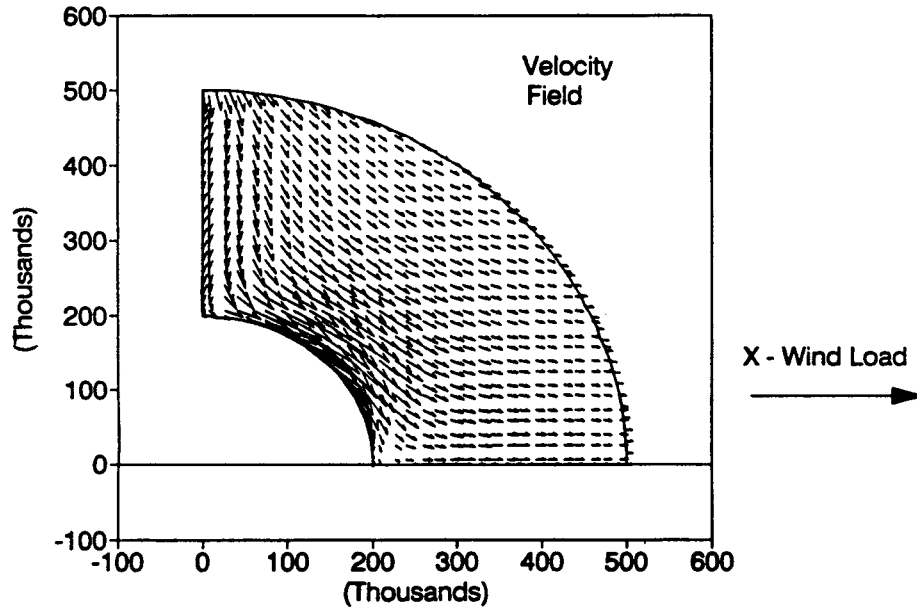


Figure 2. Wind circulation, numerical solution; quadratic bathymetry  $h = H_0 r^2$



Table 1. Norms

Iteration	Residual	Galerkin	Exact	Gradient
0	$6.515 \times 10^{-13}$	$7.225 \times 10^4$	$4.093 \times 10^{-6}$	$1.24 \times 10^7$
1	$1.554 \times 10^{-12}$	$5.548 \times 10^4$	$2.760 \times 10^{-6}$	$1.43 \times 10^5$
2	$3.586 \times 10^{-12}$	$4.233 \times 10^4$	$1.646 \times 10^{-6}$	$2.95 \times 10^4$
3	$6.666 \times 10^{-12}$	$3.315 \times 10^4$	$9.745 \times 10^{-7}$	$1.01 \times 10^4$
4	$1.254 \times 10^{-11}$	$2.624 \times 10^4$	$5.950 \times 10^{-7}$	$2.92 \times 10^3$
5	$2.578 \times 10^{-11}$	$2.027 \times 10^4$	$3.708 \times 10^{-7}$	$6.97 \times 10^2$
6	$4.691 \times 10^{-11}$	$1.582 \times 10^4$	$2.558 \times 10^{-7}$	$1.85 \times 10^2$
7	$7.186 \times 10^{-11}$	$1.291 \times 10^4$	$2.072 \times 10^{-7}$	$5.56 \times 10^1$
8	$9.896 \times 10^{-11}$	$1.102 \times 10^4$	$1.892 \times 10^{-7}$	$1.72 \times 10^1$
9	$1.279 \times 10^{-10}$	$9.778 \times 10^3$	$1.821 \times 10^{-7}$	$5.17 \times 10^0$

Table 2. Weights

Iteration	$W_1$	$W_2$	$W_3$	$W_4$	$W_5$
0	1.0	1.0	1.0	1.0	1
1	16.4	16.4	16.5	15.1	1
2	41.7	42.3	43.6	19.6	1
3	86.4	88.8	95.4	23.7	1
4	179.0	184.9	208.3	34.1	1
5	382.6	393.8	465.7	54.8	1
6	743.4	753.7	945.1	80.6	1
7	1307.1	1296.7	1744.2	107.5	1
8	2207.4	2141.2	3100.7	140.4	1
9	3722.5	3545.6	5502.8	187.9	1

solution is shown in Figure 2. A complete description of this test problem is given in Reference 16 as well as the analytic solution adapted from Reference 20.

There are five weights in this problem.  $W_1$  and  $W_2$  are the weights for the  $x$ - and  $y$ -momentum equations respectively.  $W_3$  is associated with the continuity equation.  $W_4$  and  $W_5$  are the weights associated with the prescribed  $\xi$  (on  $r = r_2$ ) and the prescribed zero normal flow on the remaining boundaries. All weights were set initially to 1.0. The weight  $W_5$  (normal flow condition) was fixed at 1.0 and the optimization process was carried out to determine  $W_1$ – $W_4$  that minimize the auxiliary function  $\Psi_a$ . Tables I and II and Figure 3 summarize the results. In Table I the numerical values of various norms are listed for each step of the optimization. The residual norm is defined as the  $L^2$ -norm of the vector  $\{\epsilon\}$  evaluated at each of the active domain collocation points and at the boundary collocation points. The Galerkin norm is actually the value of  $\Psi_a$  as defined by equation (35). The values in the exact column are the  $L^2$ -norm of the total error. This is simply the sum of the squares of the difference between the numerical and analytical solutions evaluated at the same points used for the residual norm. The gradient norm is the  $L^2$ -norm of  $\{\partial\Psi_a/\partial s_i\}$  as defined by the left-hand side of equation (42). The four norms are also plotted in Figures 3(a) and 3(b). Owing to their varying magnitudes, the values are normalized with respect to the values in the first row of Table I. The corresponding weights at each iteration are listed in Table II and plotted in Figure 3(c). For this test problem the iteration procedure was carried out until the gradient norm was less than  $10^{-6}$  of the initial gradient norm.

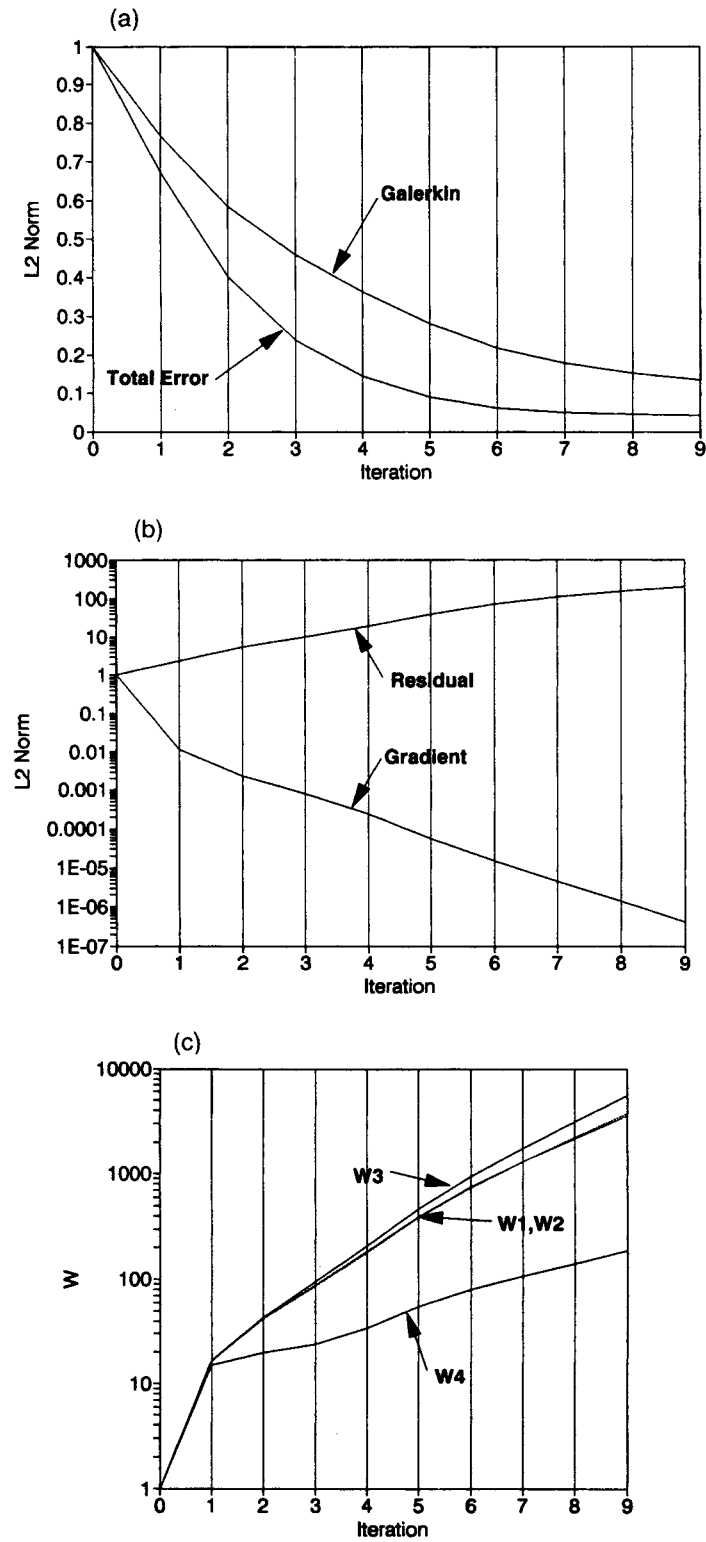


Figure 3. (a) Galerkin norm and total error norm. (b) Least squares residual norm and gradient norm. (c) Weights

## 6.2. Discussion

The significance of these results is best illustrated by Figure 3(a). We see that weights that reduce the Galerkin norm also reduce the exact (total error) norm. It is thus apparent that in the absence of any knowledge of an analytic solution we can rely on the Galerkin norm as a guide to decreasing the total error of the solution. The importance of this result is further appreciated by considering Figure 3(b), where it can be seen that the  $L^2$  residual norm of the least squares method is actually increasing as the solution becomes more accurate. Figure 3(c) reveals that although the total error can be reduced by minimizing the Galerkin norm, the weights continue to increase ( $\{\Delta s_i\}$  becomes constant) while the gradient  $\{\partial\Phi_a/\partial s_i\}$  approaches zero. This implies that the second derivatives are zero and that we are approaching a plateau of neutral stability. The selection of weights larger than those that mark the onset of this plateau has a negligible affect on the numerical solution. At the extreme condition of very large weights, however, it is conceivable that ill conditioning may cause the solution to begin to degenerate.

## 7. CONCLUSIONS

An alternative least squares method for the solution of a first-order system of partial differential equations has been presented. The new feature of the method is a systematic methodology for the determination of the optimal weights appearing in the weighted least squares method.

The results obtained for the solution of the two-dimensional shallow water equations show a superiority of the new approach when compared with the classical least squares and Galerkin methods.

## APPENDIX: PROOF OF THEOREM 2

Symbols and notations introduced in Section 4 are used in this appendix.

The proof consists of a comparison between two minimization problems. The first optimization problem does not involve weights and may be formulated as follows. Find the vector solution  $\mathbf{a}_p$  of dimension  $pN^D$  of the minimization problem

$$\begin{aligned} & \text{Min}_{\mathbf{a}} \int_{\Omega} [R_L(\mathbf{a}, \mathbf{x})]^2 d\Omega \\ & \text{subject to } \mathbf{L}_i(\mathbf{x})\mathbf{a} - f_i = 0, \quad i = 2, \dots, P, \\ & \text{subject to } \mathbf{B}_l(\mathbf{x})\mathbf{a} - g_l = 0, \quad l = 1, \dots, pb \\ & \quad \quad \quad (\text{problem } Q_h), \end{aligned} \quad (70)$$

where  $R_L(\mathbf{a}, \mathbf{x})$  is defined as the residual of the first equation:

$$R_L(\mathbf{a}, \mathbf{x}) = \mathbf{L}_1(\mathbf{x})\mathbf{a} - f_1. \quad (71)$$

The second minimization problems corresponds to the penalty formulation of problem  $Q_h$  and may be formulated as

$$\begin{aligned} & \text{Min}_{\mathbf{a}} \int_{\Omega} [R_L(\mathbf{a}, \mathbf{x})]^2 d\Omega + \sum_{i=2}^P \frac{1}{\epsilon_i} \int_{\Omega} (\mathbf{L}_i(\mathbf{x})\mathbf{a} - f_i)^2 d\mu + \sum_{l=1}^{pb} \frac{1}{v_l} \int_{\Gamma} (\mathbf{B}_l(\mathbf{x})\mathbf{a} - g_l)^2 d\mu \\ & \quad \quad \quad (\text{problem } R_s). \end{aligned} \quad (72)$$

For the continuous least squares formulation  $\mu$  represents the measure associated with the volume integral on  $\Omega$ .

The proof of Theorem 2 requires the introduction of the Lagrangian multipliers  $p_{\epsilon_i}$ ,  $p_{v_j}$ ,  $p_j$  and finally  $pb_l$ . This is the purpose of Lemmas 1 and 2; then Theorem 2 will result from the relationship between these Lagrangian multipliers.

*Lemma 1*

Problem  $Q_h$  admits a unique solution  $\mathbf{a}_p$  characterized by the existence of Lagrangian multipliers  $p_j$ ,  $j = 2, \dots, p$ , and  $pb_l = 1, \dots, pb$ , such that

$$\begin{aligned} \mathbf{A}_1 \mathbf{a}_p + \sum_{j=2}^p \frac{1}{2} \int_{\Omega} \mathbf{L}_j^T(\mathbf{x}) p_j(\mathbf{x}) d\mu + \sum_{l=1}^{pb} \frac{1}{2} \int_{\Omega} \mathbf{B}_l^T(\mathbf{x}) pb_l(\mathbf{x}) d\mu = \mathbf{f}_1, \\ \mathbf{L}_i(\mathbf{x}) \mathbf{a} - f_i = 0, \quad i = 2, \dots, p, \quad \mathbf{B}_l(\mathbf{x}) \mathbf{a} - g_l = 0, \quad l = 1, \dots, pb. \end{aligned} \quad (73)$$

*Proof.* The interior residual appearing in problem  $Q_h$  may be written in terms of the scalar product on  $\mathbb{R}^{pD}$  as

$$[R_L(\mathbf{a}, \mathbf{x})]^2 = (\mathbf{L}_1^T(\mathbf{x}) \mathbf{L}_1(\mathbf{x}) \mathbf{a}, \mathbf{a}) - 2(\mathbf{L}_1^T(\mathbf{x}) \mathbf{f}_1(\mathbf{x}), \mathbf{a}) + f_1^2(\mathbf{x}). \quad (74)$$

The last term  $f_1^2(\mathbf{x})$  does not depend on  $\mathbf{a}$ , so problem  $Q_h$  is equivalent to the minimization problem

$$\begin{aligned} \text{Min}_{\mathbf{a}} \int_{\Omega} (\mathbf{L}_1^T(\mathbf{x}) \mathbf{L}_1(\mathbf{x}) \mathbf{a}, \mathbf{a}) d\mathbf{x} - 2 \int_{\Omega} (\mathbf{L}_1^T \mathbf{f}_1, \mathbf{a}) d\mathbf{x} \\ \text{subject to } \mathbf{L}_i(\mathbf{x}) \mathbf{a} - f_i = 0, \quad \mathbf{x} \in \Omega, \quad i = 2, \dots, p, \\ \text{subject to } \mathbf{B}_l(\mathbf{x}) \mathbf{a} - g_l = 0, \quad \mathbf{x} \in \Gamma, \quad l = 1, \dots, pb, \\ (\text{problem } Q'_h). \end{aligned} \quad (75)$$

The quadratic minimization problem  $Q'_h$  admits a unique solution  $\alpha_p$  characterized by the existence of multipliers  $p_j(\mathbf{x})$ ,  $j = 2, \dots, p$ , and  $pb_l$ ,  $l = 1, \dots, pb$ , such that equation (73) is satisfied.  $\square$

*Lemma 2*

Problem  $R_h$  admits a unique solution  $\bar{\mathbf{a}}_{\epsilon}$  characterized by the existence of the Lagrangian multipliers  $p_{\epsilon_j}$  such that

$$\begin{aligned} \mathbf{A}_1 \bar{\mathbf{a}}_{\epsilon} + \sum_{j=2}^p \frac{1}{2} \int_{\Omega} \mathbf{L}_j(\mathbf{x}) \mathbf{p}(\mathbf{x})_{\epsilon_j} d\mu + \sum_{l=1}^{pb} \frac{1}{2} \int_{\Gamma} \mathbf{B}_l(\mathbf{x}) \mathbf{p}(\mathbf{x})_{v_l} d\mu = \mathbf{f}_1, \\ \mathbf{p}(\mathbf{x})_{\epsilon_i} = \frac{2}{\epsilon_i} (\mathbf{L}_i(\mathbf{x}) \mathbf{a} - f_i), \quad i = 2, \dots, p, \quad \mathbf{p}(\mathbf{x})_{v_l} = \frac{2}{v_l} (\mathbf{B}_l(\mathbf{x}) \mathbf{a} - g_l), \quad l = 1, \dots, pb. \end{aligned} \quad (76)$$

*Proof.* Define the objective function of problem  $R_h$  as

$$\begin{aligned} f(\mathbf{a}, \epsilon_1, \dots, \epsilon_p, v_1, \dots, v_{pb}) = \int_{\Omega} [R_L(\mathbf{a}, \mathbf{x})]^2 d\Omega + \sum_{i=2}^p \frac{1}{\epsilon_i} \int_{\Omega} (\mathbf{L}_i(\mathbf{x}) \mathbf{a} - f_i)^2 d\mu \\ + \sum_{l=1}^{pb} \frac{1}{v_l} \int_{\Gamma} (\mathbf{B}_l(\mathbf{x}) \mathbf{a} - g_l)^2 d\mu. \end{aligned} \quad (77)$$

The function  $\mathbf{a} \rightarrow f(\mathbf{a}, \epsilon_1, \epsilon_2, \dots, \epsilon_p, \nu_1, \dots, \nu_{pb})$  defined from  $\mathbb{R}^{(p(D+1)+pb)} \rightarrow \mathbb{R}$  is a differentiable function, so the first-order condition for a minimum corresponds to a null differential at the solution  $\bar{\mathbf{a}}$ .

After differentiation one can easily get

$$\mathbf{A}_1 \bar{\mathbf{a}}_\epsilon + \sum_{j=2}^p \frac{1}{\epsilon_j} \int_{\Omega} \mathbf{L}_j^T(\mathbf{x})(\mathbf{L}_j(\mathbf{x})\bar{\mathbf{a}}_\epsilon - \mathbf{f}_j) d\mu + \sum_{l=1}^{pb} \frac{1}{\nu_l} \int_{\Gamma} \mathbf{B}_l^T(\mathbf{x})(\mathbf{B}_l(\mathbf{x})\bar{\mathbf{a}}_\epsilon - \mathbf{g}_l) d\Gamma = \mathbf{f}_1. \quad (78)$$

By defining  $p_{\epsilon_j}$  as in (76), equations (76) may be easily found from equation (78).  $\square$

By virtue of the implicit function theorem, it may be shown that the vectors  $p_\epsilon(\mathbf{x}, c)$  and the  $\bar{\mathbf{a}}_\epsilon$  are analytic functions of the variables  $\epsilon_j, j = 2, \dots, p$ . This last property permits one to prove the following lemma.

*Lemma 3*

Let  $\bar{\mathbf{a}}_\epsilon$  and  $p_{\epsilon_j}$  be the solution vectors associated with equations (76) and let  $\mathbf{a}_p$  and  $p_j$  be the solution of equations (73). Then, if  $p_{\epsilon_j}$  depends on  $\epsilon_j$  only ( $j = 2, \dots, p$ ) and  $pb_{\nu_j}$  depends on  $\nu_j$  only ( $j = 1, \dots, pb$ ),

$$\begin{aligned} p_{\epsilon_j}(\mathbf{x}) - p_j(\mathbf{x}) &= H_j(\mathbf{x})\epsilon_j + O(\epsilon_j^2), \quad j = 2, \dots, p, \\ pb_{\nu_l}(\mathbf{x}) - pb_l(\mathbf{x}) &= D_l(\mathbf{x})\nu_l + O(\nu_l^2), \quad l = 1, \dots, pb, \end{aligned} \quad (79)$$

where  $D_l(\mathbf{x})$  and  $H_j(\mathbf{x})$  are functions independent of the penalty parameters  $\nu_l$  and  $\epsilon_j$  respectively.

*Proof*

$$\begin{aligned} p_{\epsilon_j}(\mathbf{x}) &= p_{0j}(\mathbf{x}) + \epsilon_j p_{1j}(\mathbf{x}) + O(\epsilon_j^2), \quad j = 2, \dots, p, \\ pb_{\nu_l}(\mathbf{x}) &= pb_{0l}(\mathbf{x}) + \nu_l pb_{1l}(\mathbf{x}) + O(\nu_l^2), \quad l = 1, \dots, pb. \end{aligned} \quad (80)$$

By virtue of the uniqueness of the solution of problem  $Q_h$ , one can easily see that  $p_{0j}(\mathbf{x})$  and  $pb_{\nu_l}(\mathbf{x})$  verify equations (73) and therefore  $p_{0j}(\mathbf{x}) = p_j(\mathbf{x})$  and  $pb_{0l}(\mathbf{x}) = pb_l(\mathbf{x})$ . Define now  $p_{1j}(\mathbf{x}) = H_j(\mathbf{x})$  and  $pb_{1l}(\mathbf{x}) = D_l(\mathbf{x})$ . Lemma 3 is proved.  $\square$

The last lemma permits us to express the Lagrangian multiplier vectors  $p_{\epsilon_j}(\mathbf{x}), j = 2, \dots, p$ , in terms of the penalty weights  $\epsilon_j$  and to express  $pb_{\nu_j}(\mathbf{x}), j = 1, \dots, pb$ , in terms of the penalty weights  $\nu_j$ :

$$\begin{aligned} p_{\epsilon_j}(\mathbf{x}) &= p_j(\mathbf{x}) + \epsilon_j H_j(\mathbf{x}) + O(\epsilon_j^2), \quad j = 1, \dots, p, \\ pb_{\nu_l}(\mathbf{x}) &= pb_l(\mathbf{x}) + \nu_l D_l(\mathbf{x}) + O(\nu_l^2), \quad l = 1, \dots, pb, \end{aligned} \quad (81)$$

By definition of  $p_{\epsilon_j}(\mathbf{x})$  and  $pb_{\nu_l}(\mathbf{x})$  (equations (76)) and after integration over  $\Omega$  one can get

$$\begin{aligned} \int_{\Omega} (\mathbf{L}_j(\mathbf{x})\bar{\mathbf{a}}_\epsilon - \mathbf{f}_j) d\Omega &= \frac{\epsilon_j}{2} \int_{\Omega} p_j(\mathbf{x}) d\Omega + O(\epsilon_j^2), \quad j = 2, \dots, p, \\ \int_{\Omega} (\mathbf{B}_l(\mathbf{x})\bar{\mathbf{a}}_\epsilon - \mathbf{f}_l) d\Omega &= \frac{\nu_l}{2} \int_{\Omega} pb_l(\mathbf{x}) d\Omega + O(\nu_l^2), \quad j = 1, \dots, p. \end{aligned} \quad (82)$$

The last equation finishes the proof of Theorem 2 for the continuous case. For the collocation case the proof is similar, only the measure  $\mu$  has to be changed in terms of the set of collocation points  $\mathbf{x}_p$ ,

$i = 1, \dots, k$ , for the interior collocation points and  $\mathbf{x}_i, i = k+1, \dots, m$ , for the boundary collocation points:

$$\int_{\Omega} f(\mathbf{x}) d\mu = \sum_{i=1}^k v_i f(\mathbf{x}_i), \quad \int_{\Gamma} f(\mathbf{x}) d\mu = \sum_{i=k+1}^m f(\mathbf{x}_i). \quad (83)$$

## REFERENCES

1. R. Glowinski, *Numerical Methods for Nonlinear Variational Problems*, Springer, Berlin, 1986.
2. J. Bensabat and D. G. Zeitoun, 'A least-squares formulation for the solution of transport problems', *Int. j. numer. methods fluids*, **10**, 623–636.
3. T. J. R. Hughes, L. P. Franca and G. M. Hulbert, 'A new finite element formulation for computational fluid dynamics: the Galerkin/least squares method for advective–diffusive equations', *Comput. Methods Appl. Mech. Eng.* **73**, 173–189 (1989).
4. E. D. Eason, 'A review of least-squares methods for solving partial differential equations', *Int. j. numer. methods eng.*, **10**, 1021–1046 (1976).
5. K. O. Friedrichs, 'Symmetric positive differential equations', *Commun. Pure Appl. Math.*, **11**, 333–418 (1958).
6. L. Lapidus and G. F. Pinder, *Numerical Solution of Partial Differential Equations in Science and Engineering*, Wiley, New York, 1982.
7. P. Lesaint and P. A. Raviart, 'Finite element collocation methods for first order systems', *Math. Comput.*, **33**, 891–918 (1979).
8. A. K. Aziz and J. L. Liu, 'A weighted least squares method for the backward–forward heat equation', *SIAM J. Numer. Anal.*, **28**, 156–167 (1991).
9. J. H. Bramble and A. H. Shatz, 'Least squares methods for  $2m$ th order elliptic boundary value problems', *Math. Comput.*, **25**, 1–32 (1971).
10. W. L. Wendland, *Elliptic Systems in the Plane*, Pitman, London, 1979.
11. A. K. Aziz, R. B. Kellogg and A. B. Stephens, 'Least squares methods for elliptic systems', *Math. Comput.*, **44**, 53\_70 (1985).
12. G. F. Carey and B. N. Jiang, 'Least-squares finite elements for first order hyperbolic systems', *Int. j. numer. methods eng.*, **26**, 81–93 (1988).
13. B. N. Jiang and G. F. Carey, 'A stable least-squares finite element method for non-linear hyperbolic problems', *Int. j. numer. methods fluids*, **8**, 933–942 (1988).
14. B. N. Jiang and G. F. Carey, 'Least-squares finite element methods for compressible Euler equations', *Int. j. numer. methods fluids*, **10**, 557–568 (1990).
15. J. P. Laible and G. F. Pinder, 'Least squares collocation solution of differential equations on irregularly shaped domains using orthogonal meshes', *Numer. Methods PDEs*, **5**, (1989).
16. J. P. Laible and G. F. Pinder, 'Solution of the shallow water equations by least-squares collocation', *Adv. Water Resources Res.*, submitted.
17. P. A. Raviart and J. M. Thomas, *Introduction a l'Analyse Numerique des Equations aux Derivees Partielles*, Masson, Paris, 1988.
18. T.-F. Chen, 'On least-squares approximations to compressible flow problems', *Numer. Methods PDEs*, **2**, 207–228 (1986).
19. D. G. Zeitoun, J. P. Laible and G. F. Pinder, 'An iterative least-squares solution for boundary value problems', *Int. j. numer. methods eng.*, submitted.
20. D. R. Lynch and W. G. Gray, 'Analytic solutions for computer flow model testing', *J. Hydraul. Div., ASCE*, **104**, 1409–1428 (1978).