

Non-oscillatory and Non-diffusive Solution of Convection Problems by the Iteratively Reweighted Least-Squares Finite Element Method

BO-NAN JIANG

ICOMP, NASA Lewis Research Center, Cleveland, Ohio 44135

Received October 16, 1990; revised October 11, 1991

In this paper we first introduce the least-squares finite element method (LSFEM) for two-dimensional steady-state pure convection problems with smooth solutions and compare the LSFEM with other finite element methods. We prove that the LSFEM has the same stability estimate as the original equation; i.e., the LSFEM has better control of the streamline derivative than the streamline upwinding Petrov-Galerkin method. Numerical convergence rates are given to show that the LSFEM is almost optimal. Then we use this LSFEM as a framework to develop an iteratively reweighted least-squares finite element method (IRLSFEM) to obtain non-oscillatory and non-diffusive solutions for problems with contact discontinuities. This new method produces a highly accurate numerical solution that has a sharp discontinuity in one element. A number of examples solved by using triangular and bilinear elements are presented to show that the method can convect contact discontinuities without error. © 1993 Academic Press, Inc.

1. INTRODUCTION

In this paper we introduce and test numerically the iteratively reweighted least-squares finite element method for the solution of two-dimensional steady-state pure convection problems. This new method is designed to obtain accurate non-oscillatory discontinuous solutions.

We shall consider the following steady-state boundary value problem,

$$u_\beta = 0 \quad \text{in } \Omega, \quad (1a)$$

$$u = g \quad \text{on } \Gamma_- \quad (1b)$$

where Ω is a bounded convex domain in \mathbb{R}^2 with boundary Γ , $u = u(x, y)$ is the dependent variable (e.g., the concentration), $\beta = (\beta_1, \beta_2)$ is a specified convection vector which may be constant or spatially varying, $u_\beta = \beta \cdot \nabla u$ denotes the derivative in β direction and g is the given data on the inflow boundary Γ_- defined by

$$\Gamma_- = \{(x, y) \in \Gamma : \mathbf{n}(x, y) \cdot \beta < 0\},$$

in which \mathbf{n} is the outward unit normal to Γ at point $(x, y) \in \Gamma$. The problem (1) is purely hyperbolic.

For convenience of discussion, at first let us assume that the convection vector β is a constant. In this case, the characteristics of the problem (1) are straight lines parallel to β , and the analytic solution of problem (1) is constant along a characteristic. The value of this constant is equal to the given value of g at the intersection of this characteristic and the inflow boundary. Thus the solution will be discontinuous with a jump across a characteristic, if the boundary data g is discontinuous. As a simple 2D testing problem, this situation still represents a great challenge to researchers of numerical shock-capturing techniques.

Commonly used numerical methods for hyperbolic problems are of the following types (see, e.g., Fletcher [10], Hirsch [13], Johnson [21], and Pironneau [31]): method of characteristics, finite difference, and finite element methods. In principle the method of characteristics is very good, but it is rather cumbersome in practice. Usually one uses finite difference or finite element methods based on a mesh which may not be adapted to fit the characteristics of the particular problem. In such a case, if the exact solution has a jump discontinuity (contact discontinuity) across a characteristic, all conventional finite difference and finite element methods will produce approximate solutions which either oscillate or smear out a sharp front. Finding accurate approximations of the discontinuous solutions of hyperbolic equations has been a persistent difficult task in modern numerical mathematics and computational physics.

One research direction towards better resolution around discontinuities is to use an adaptive h -refinement or remeshing strategy, such as that extensively investigated by Oden and Demkowicz [29] and Peraire *et al.* [30]. However, the data structure and the programming of h -refinement are complicated, especially for three-dimensional problems.

Impressive resolution may be obtained by using filter methods [4, 23, 9]. For example, one may use any good

finite difference schemes, such as the schemes developed by Deconinck *et al.* [7], to obtain an initial approximation and then use filter methods to sharpen the discontinuities.

Another potential way is to use conventional finite difference or finite element methods to obtain approximate solutions and then apply imaging processing techniques to detect the locations of discontinuities [32, 33].

All of the above-mentioned methods can only improve the resolution around the contact discontinuities. The L_1 procedure for non-oscillatory solutions, first proposed by Lavery in [24, 25] for one-dimensional problems, is a rather different approach. The L_1 solutions are highly accurate and right up to the edge of the discontinuity. Lavery uses linear programming to minimize the L_1 norm of the residuals of the overdetermined algebraic equations. Standard finite difference leads to determined linear algebraic systems. In order to obtain an overdetermined system Lavery relies on non-traditional tricks, such as gradually adding a small viscous term to one-dimensional Burgers' equation [24]. However, it is very difficult to extend these tricks to two-dimensional problems. Another difficulty in Lavery's L_1 procedure is that the linear programming algorithm of Barrrodale and Roberts [1] is very expensive. This excludes the possibility of practical use of the L_1 method.

Because of the difficulties associated with the L_1 procedure in both theory and calculation, we completely avoid the application of the L_1 concept in this study. Our procedure for obtaining non-oscillatory solutions is based on the least-squares method. The idea can be explained as follows. In the usual least-squares curve fitting, the least-squares procedure does its best in a sense of least-squares of the residual to make the curve pass through or by all of the data. If the data are smooth, the least-squares fitting yields a very good approximation. However, if the data contain abrupt changes, the least-squares procedure will produce an oscillatory and diffusive curve around sharp changes. In such a case, the trouble comes from the fact the least-squares fitting makes the use of individual datum equally important. If we give up the outliers in the data and require the remaining data be satisfied exactly, the curve will be more smooth. The same thing happens in the least-squares solutions of discretized hyperbolic equations. The least-squares method treats equally the equations in "shocked" elements (in which the discretized scheme does not hold) and "smooth" elements. If we can identify the "shocked" elements and permit the equations in the "shocked" elements not to be satisfied while requiring that the remaining equations be satisfied exactly, then this modified least-squares solution will not have oscillation.

In this paper we first introduce the least-squares finite element method for two-dimensional steady-state pure convection problems with smooth solutions and compare it with other finite element methods. We prove that the

least-squares method has the same stability estimate as the original equation; that is, the least-squares method has better control of the streamline derivative than the streamline upwinding Petrov-Galerkin method. Numerical convergence rates are given to show that the least-squares method is almost optimal.

Since the LSFEM produces a very good initial approximation to the exact solution, we use this information to find "shocked" elements. Then we use the least-squares method again. But this time we put a small weight for "shocked" elements to eliminate their "pollution," and repeat this procedure a few times until a convergent discontinuous solution is reached. This is our iteratively reweighted least-squares finite element method.

The reweighting must be based on an overdetermined system. Fortunately, it is trivial to have an overdetermined system in the least-squares finite element method. This can be explained as follows. The LSFEM with numerical quadrature is equivalent to a weighted collocation least-squares method [3], in which at first the residual equations are collocated at the interior points in each element and then the algebraic system is approximately solved by the weighted least-squares method. The Gaussian points for calculating the element matrices in the LSFEM correspond to the collocation points in collocation methods. If the order of Gaussian quadrature (or the number of quadrature points) is appropriately chosen, the least-squares finite element method amounts to solving an overdetermined system.

The arrangement of this paper is as follows. The least-squares method, its comparison with other methods, and the convergence tests for smooth problems are presented in Section 2. In Section 3 we describe the iteratively reweighted least-squares finite element method. The numerical results in Section 4 contain non-oscillatory solutions for problems with constant or spatially varying convection vector fields on uniform or unstructured meshes. In Section 5 we discuss the limitation of the method in the present implementation and the possibility of further improvement. Conclusions are drawn in Section 6.

2. THE LEAST-SQUARES FINITE ELEMENT METHOD

2.1. Preliminaries and Notations

As we have already shown in our previous papers (see [20] and the references therein), the LSFEM is a universal method for numerical solution of partial differential equations. It does not matter whether the partial differential equations are elliptic, parabolic, or hyperbolic. As long as the partial differential equation has a unique solution, the LSFEM always gives a reasonably good approximate solution. The work done in this paper is a natural extension of our least-squares method.

The problem (1) can be taken as a time-dependent problem, if we consider one space coordinate as a time-like coordinate. Then we may use the implicit time-marching least-squares finite element method introduced in [2] to obtain an approximate solution. The time marching is necessary for the Euler equations in aerodynamics [18, 19], since the Euler equations have nonunique solutions. The time-marching least-squares finite element method implicitly introduces an artificial dissipation to exclude the solutions with expansion shocks. For the linear hyperbolic problem (1), the time-marching is not necessary, because it has a unique solution, continuous or discontinuous. For this reason, we rather treat the problem (1) as two-dimensional and directly apply the LSFEM to attack it. The general formulation of least-squares finite element methods for first-order partial differential equations can be found in [20].

In order to see the advantages of the LSFEM and compare it with other finite element methods, let us consider the following linear hyperbolic equation:

$$u_\beta + u = f \quad \text{in } \Omega, \quad (2a)$$

$$u = g \quad \text{on } \Gamma_-, \quad (2b)$$

where f is a given source function. Without loss of generality we assume that the boundary data g is zero.

Throughout this paper, we use the following notations: $L_2(\Omega)$ denotes the space of square-integrable functions defined on Ω with the inner product

$$(u, v) = \int_{\Omega} uv \, d\Omega \quad u, v \in L_2(\Omega),$$

and the norm

$$\|u\|^2 = (u, u), \quad u \in L_2(\Omega).$$

$H^r(\Omega)$ denotes the Sobolev space of functions with square-integrable derivatives of order up to r ; $\|\cdot\|_r$ denotes the usual norm for $H^r(\Omega)$. We also use the following notations:

$$\begin{aligned} \langle u, w \rangle &= \int_{\Gamma} uw \mathbf{n} \cdot \beta \, ds, \\ \langle u, w \rangle_+ &= \int_{\Gamma_+} uw \mathbf{n} \cdot \beta \, ds, \\ |u|_{\Gamma} &= \left(\int_{\Gamma} u^2 |\mathbf{n} \cdot \beta| \, ds \right)^{1/2}, \end{aligned}$$

where

$$\Gamma_+ = \{(x, y) \in \Gamma : \mathbf{n}(x, y) \cdot \beta \geq 0\}.$$

We note that by Green's formula

$$(u_\beta, w) = \langle u, w \rangle - (u, w_\beta).$$

Further, we also define the function space

$$S = \{u \in H^1(\Omega) : u = 0 \text{ on } \Gamma_-\},$$

and the corresponding finite element subspace S_h ; i.e., S_h is the space of continuous piecewise polynomial functions of order k . Here the parameter h represents the maximal diameter of the elements. By the finite element interpolation theory [6, 28] we have: Given a function $u \in H^{k+1}(\Omega)$, there exists an interpolant $\hat{u}^h \in S_h$ such that

$$\|u - \hat{u}^h\| \leq ch^{k+1} \|u\|_{k+1}, \quad (3a)$$

$$\|\nabla u - \nabla \hat{u}^h\| \leq ch^k \|u\|_{k+1}; \quad (3b)$$

here and below c denotes a constant independent of the mesh parameter h , with possibly different values in each appearance.

2.2. The Standard Galerkin Method

Now let us look at the following standard Galerkin method for the problem (2) (see, e.g., Johnson [21]). Find $u^h \in S_h$ such that

$$(u_\beta^h + u^h, w^h) = (f, w^h) \quad \forall w^h \in S_h. \quad (4)$$

We define the error $e = u - u^h$. Then the error estimate for the standard Galerkin method is

$$\|e\| + |e|_{\Gamma} \leq ch^k \|u\|_{k+1}, \quad (5)$$

which is one order lower than that for elliptic and parabolic problems. Furthermore, in the continuous problem (2), we have the following stability estimate:

$$\|u\| + \|u_\beta\| + |u|_{\Gamma} \leq c \|f\|. \quad (6)$$

But in the standard Galerkin method the stability estimate is

$$\|u^h\| + |u^h|_{\Gamma} \leq c \|f\|, \quad (7)$$

which has no control of $\|u_\beta^h\|$.

2.3. The SUPG Method

In order to obtain better accuracy and stability, methods of upwinding type have appeared, see, e.g., papers by Dendy [8], Wahlbin [34], Christie *et al.* [5], Hughes and Brooks [15], Johnson *et al.* [22], Morton and Parrot [27]. Below

we shall look at the streamline upwinding Petrov–Galerkin method (SUPG) (Hughes [14]) or the streamline diffusion method (Johnson [21]): Find $u^h \in S_h$ such that

$$(u_\beta^h + u^h, w^h + hw_\beta^h) = (f, w^h + hw_\beta^h) \quad \forall w^h \in S_h. \quad (8)$$

Johnson and his colleagues have derived the error estimate,

$$\left(\|e\|^2 + h \|e_\beta\|^2 + \frac{(1+h)}{2} |e|_T^2 \right)^{1/2} \leq ch^{k+1/2} \|u\|_{k+1}, \quad (9)$$

which is near optimal. However, in the SUPG the corresponding stability estimate is

$$\|u^h\| + \sqrt{h} \|u_\beta^h\| + |u^h|_T \leq c \|f\|, \quad (10)$$

which means that the streamline derivative is less controlled than in the original problem. Another disadvantage of the SUPG in practical calculation is that the stiffness matrix is non-symmetric, which makes the solution of large-scale problems very difficult.

2.4. The Least-Squares Method

Now let us introduce the LSFEM. We assume that $f \in L_2(\Omega)$. For an arbitrary trial function $v \in S$, we define the residual function $R = v_\beta + v - f$. The least-squares method is based on minimizing the residual function in a least-squares sense. We construct the following quadratic functional:

$$\begin{aligned} I(v) &= \|R\|^2 = \|v_\beta + v - f\|^2 \\ &= (v_\beta + v - f, v_\beta + v - f). \end{aligned} \quad (11)$$

The least-squares method reads: Find $u \in S$ such that

$$I(u) \leq I(v) \quad \forall v \in S.$$

Taking variation of I with respect to v and setting $\delta I = 0$ and $\delta v = w$, lead to the least-squares weak statement: Find $u \in S$ such that

$$b(u, w) = l(w) \quad \forall w \in S, \quad (12)$$

where $b(u, w) = (u_\beta + u, w_\beta + w)$ and $l(w) = (f, w_\beta + w)$. The corresponding LSFEM has the following form: Find $u^h \in S_h$ such that

$$b(u^h, w^h) = l(w^h) \quad \forall w^h \in S_h. \quad (13)$$

Let us now turn to the error estimate. Since we can replace w in (12) by w^h , we have

$$b(u, w^h) = l(w^h) \quad \forall w^h \in S_h. \quad (14)$$

By subtracting (13) from (14) we obtain the following orthogonality for the error e :

$$b(e, w^h) = 0.$$

Let $\hat{u} \in S_h$ be the interpolant of u satisfying (3) and write $\rho = u - \hat{u}^h$ and $\theta = \hat{u}^h - u^h$ so that $e = \rho + \theta$. Then we have

$$\begin{aligned} \|e_\beta + e\|^2 &= b(e, e) = b(e, \rho) + b(e, \theta) = b(e, \rho) \\ &\leq \|e_\beta + e\| \|\rho_\beta + \rho\| \end{aligned}$$

or

$$\|e_\beta + e\| \leq \|\rho_\beta + \rho\| \leq \|\rho_\beta\| + \|\rho\|.$$

Recalling (3) we obtain the error estimate:

$$\|e_\beta + e\| \leq ch^k \|u\|_{k+1}. \quad (15)$$

Since the residual of approximate solution $R^h = u_\beta^h + u^h - f = e_\beta + e$, (15) is also the residual estimate,

$$\|R^h\| \leq ch^k \|u\|_{k+1}, \quad (16)$$

which means that the residual estimate is optimal. Using Green's formula and the boundary condition, (16) can be rewritten as

$$(\|e\|^2 + \|e_\beta\|^2 + \langle e, e \rangle_+)^{1/2} \leq ch^k \|u\|_{k+1}, \quad (17)$$

which shows that the error estimate for e_β is optimal, but the error estimate for e is one order lower than optimal. Although in numerical tests (see below) we have observed that the accuracy of the least-squares method is higher than the k th order, it is still an open question for getting a better theoretical error estimate, in general.

By taking $w^h = u^h$ in (13), we can obtain the stability estimate:

$$\|u^h\| + \|u_\beta^h\| + |u^h|_T \leq c \|f\|. \quad (18)$$

This estimate is the same as the above estimate (6) for the original problem (2). It means that the least-squares method has better control of the streamline derivative than the SUPG. We also note that the bilinear form in (13) is symmetric. Therefore, the matrix of the resulting algebraic system is symmetric and positive definite. This is a very important advantage of the least-squares method over other methods in practice.

2.5. Numerical Experiments of the Least-Squares Method

We chose the model problem,

$$\frac{\partial u}{\partial x} + \frac{\partial u}{\partial y} = \sin(x + y) \quad \text{in } \Omega, \quad (19a)$$

$$u = 0 \quad \text{on } \Gamma_-, \quad (19b)$$

where $\Omega = \{(x, y) \in \mathbb{R}^2: 0 < x < 1, 0 < y < 1\}$ is the unit square, and $\Gamma_- = \{(x, y) \in \Gamma: x = 0 \text{ or } y = 0\}$, in which Γ is the boundary. This problem has a smooth exact solution $u = \sin(x) \sin(y)$.

We have tested the bilinear element on uniform meshes. At first the one-point Gaussian quadrature was used for calculating stiffness matrices. In this case, the least-squares method is equivalent to the collocation least-squares method with one collocation point at the center of each element. It is easy to check that in such a collocation method, the number of discretized algebraic equations “*Nequ*” is equal to the number of unknowns “*Nelem*,” here “*Nelem*” is the number of elements. In other words, in such a case we solve a determined system. Therefore, there is no difference between the least-squares solution and the direct collocation solution. The numerical result for convergence rate is shown in Fig. 1. The optimal rate, i.e., $\|e\| \leq ch^2$, is observed.

Also we would like to mention that the least-squares method with one-point quadrature is equivalent to the central finite difference scheme. Therefore, the optimality of the LSFEM may be derived by using the finite difference theory.

The numerical rate of convergence with the 2×2 Gauss rule is also included in Fig. 1. In this case, the LSFEM solves an overdetermined system. The convergence rate is around $O(h^{1.75})$, which is near optimal. Here, more theoretical study is needed.

We also did the numerical tests with specified extra boundary conditions on the outflow boundary $\Gamma_+ = \{(x, y) \in \Gamma: x = 1 \text{ or } y = 1\}$. In this case, the least-squares method with the 2×2 Gauss rule gives the optimal rate of convergence $\|e\| \leq ch^2$ (Fig. 1).

From numerical experiments we may conclude that the LSFEM for the hyperbolic equation has an optimal or near optimal rate of convergence depending on the number of Gaussian points in calculation of element matrices. More

Gaussian points yield slightly less optimal results, because the least-squares method is not able to make the residual of each equation in the overdetermined system equal to zero.

3. THE ITERATIVE REWEIGHTED LEAST-SQUARES FINITE ELEMENT METHOD

If the exact solution is discontinuous, the above least-squares method still performs quite well. Of course, the argument about the error estimates does not hold. As expected, the least-squares method smears out the jump across a characteristic. As we pointed out in Section 1, the trouble comes from “shocked” elements, where the direction derivative across the jump approaches infinity and the discretized equation is not valid. But the usual least-squares method does not recognize them and just treats equally “shocked” and “smooth” elements. This is the reason that we would like to use the reweighting to suppress the interference of “shocked” elements.

Let us consider the problem (1). Our iteratively reweighted least-squares method is based on repeatedly solving a weighted least-squares problem: Find the minimizer \bar{u}^h of

$$I(u^h) = \sum_{j=1}^{N_{elem}} \left(\sum_{l=1}^{N_{gaus}} W_l w_l (R_l)^2 |J(\xi_l, \eta_l)| \right)_j, \quad (20)$$

in which

$$W_l = \frac{1}{|R_l|_{previous}^6} \quad (21)$$

and

$$R_l = u^h(\xi_l, \eta_l) - f(\xi_l, \eta_l), \quad (22)$$

where W_l denotes the assigned weight, and R_l stands for the residual at each Gaussian point, N_{gaus} denotes the number of Gaussian points, w_l is the weight of Gaussian quadrature, $|J|$ is the determinant of the Jacobian matrix, and (ξ_l, η_l) are the local coordinates of Gaussian points. As usual, u^h can be expressed as

$$u^h(\xi, \eta) = \sum_{m=1}^{N_{node}} \Psi_m(\xi, \eta) U_m, \quad (23)$$

where “ N_{node} ” is the number of nodes in an element, Ψ_m denotes the shape function, and U_m is the nodal value. In order to make the problem (20) meaningful, we must have an overdetermined algebraic system. It can be realized simply by appropriately choosing the number of Gaussian points.

The IRLSFEM would begin with the initial weight

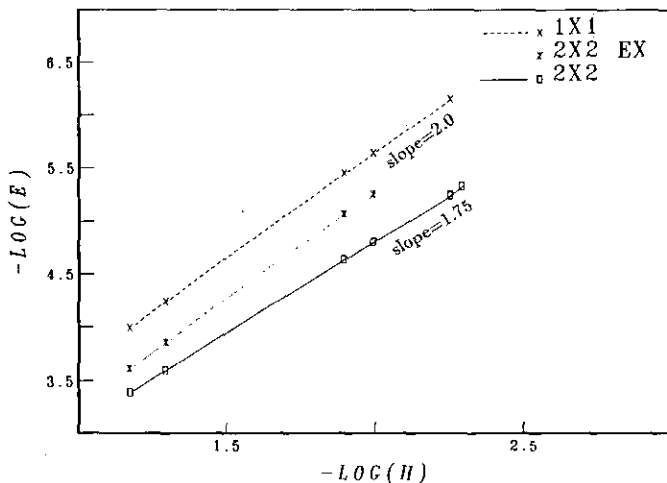


FIG. 1. Computed convergence rate for the pure convection problem.

$W_i=1$. This first step is nothing but the least-squares method introduced in Section 2. The result then determines a new set of weights by (21). In the second iteration, the residual $|R_i|$ is larger in "shocked" elements. Thus the weight W_i for "shocked" elements is smaller, and their inference becomes less important. This procedure is repeated until $\|\bar{u}_{\text{current}}^h - \bar{u}_{\text{previous}}^h\|$ is small. Our numerical experiments reveal that the difference between the residuals of "shocked" elements and those of their neighboring elements in the first least-squares solution is not significant enough. We put the sixth power in (21) in order to additionally increase the importance of "smooth" elements and reduce the contamination of "shocked" elements, and also accelerate the convergence. This is reasonable, since we want to completely eliminate the equations, which will have nonzero residuals, from the system. This trick is applicable, also because our non-weighted least-squares method is good enough to locate the "shocked" elements. That is, in the results of our least-squares method, the absolute value of the residuals in "shocked" elements is always greater than that in other elements.

We may further simplify the procedure by using another simple and reliable "shock" indicator—the variation of nodal values in each element—instead of the residual. The variation is defined as

$$V = \sum_{m=1}^{N_{\text{node}}} |U_m - U_{m-1}|, \quad U_0 = U_{N_{\text{node}}}. \quad (24)$$

Therefore, the following weight is suggested:

$$W_i = \begin{cases} 10^{15} & \text{if } |V|_{\text{previous}}^6 < 10^{-7}, \\ \frac{1}{|V|_{\text{previous}}^6} & \text{otherwise.} \end{cases} \quad (25)$$

Here some measures have been taken to prevent overflow. The advantage of using the variation as a "shock" indicator is as follows: Once the jump in the boundary data is given, we may know the exact values of the variation in "shocked" elements in advance. There are only a few possible values which depend on the type of finite element and are independent of the shape and size of the particular element and have no relation with the location of quadrature points.

The implementation of this reweighted least-squares method is really straightforward. If a least-squares finite element code is already available, it needs only a few additional line of FORTRAN statements.

4. NUMERICAL RESULTS

4.1. Constant Convection Field

As the first example, we consider the following problem with a constant convection vector,

$$\frac{\partial u}{\partial x} + (\tan 35^\circ) \frac{\partial u}{\partial y} = 0 \quad \text{in } \Omega, \quad (26a)$$

where $\Omega = \{(x, y) \in \mathbb{R}^2: 0 < x < 1, 0 < y < 1\}$ is the unit square with the boundary Γ . The inflow boundary conditions are

$$u = 2 \quad \text{on } \Gamma_1 = \{(x, y) \in \Gamma: x = 0\}, \quad (26b)$$

$$u = 1 \quad \text{on } \Gamma_2 = \{(x, y) \in \Gamma: x > 0 \text{ and } y = 0\}, \quad (26c)$$

Equations (26) represent uniform flow along straight lines inclined at an angle of 35° with respect to the x -axis. The jump discontinuity occurs along the line $y = x \tan 35^\circ$. In this case, the analytic solution is

$$u = 2 \quad \text{on and above the line } y = x \tan 35^\circ,$$

$$u = 1 \quad \text{below the line } y = x \tan 35^\circ.$$

The boundary conditions (26b) and (26c) can be transferred into the source term in Eq. (26a). Therefore, the formulation of the least-squares method described in Section 2 can be applied to the problem (26).

Most of the computational results presented in this paper were obtained in double precision from a PC-386. A direct solver with variable band-width was used to obtain the solution of linear algebraic equations. The computing time will be significantly shortened by using the preconditioned conjugate gradient method [17], since the least-squares solution is already close to the accurate solution and final iterations are often just for correcting one or two nodal values which have not yet reached 15-digit accuracy.

4.1.1. *Linear triangular element.* Numerical experiments were carried out for the problem (26) using linear triangular elements on uniform meshes with $n = 5, 15$. Here n is the number of grids in each coordinate. For triangular elements, we use the one-point Gaussian quadrature. There are $2n^2$ elements, so we have $2n^2$ equations. Since there are $(n+1)^2$ nodal values and $(2n+1)$ boundary conditions, the number of unknowns is $(n+1)^2 - (2n+1) = n^2$. That is, the number of equations is double the number of unknowns. Therefore, the least-squares method amounts to solving an overdetermined system. It does not make sense to take more quadrature points, because in a linear triangular element $\partial u^h / \partial x$ and $\partial u^h / \partial y$ are constants, and thus the residuals at different points are the same.

The least-squares results for $n = 5$ (50 triangles) are listed in Table I. The numbers in Table I are the nodal values. Because the mesh is very coarse, the jump discontinuity is smeared severely. Starting from this bad least-squares solution, after four iterations of the IRLSFEM, we obtained the perfect discontinuous solution listed in Table II. This solution is correct to 15 digits. Here we should note that a

TABLE I
Nodal Values of LSFEM Solution for Constant Convection Problem (50 Triangles)

Table with 6 columns and 6 rows of numerical data representing nodal values for a 50-triangle problem.

TABLE II
Nodal Values of IRLSFEM Solution for Constant Convection Problem (50 Triangles)

Table with 6 columns and 6 rows of numerical data, mostly showing values of 2.0000000000000000.

TABLE III
Nodal Values of LSFEM Solution for Constant Convection Problem (450 Triangles)

Table with 15 columns and 15 rows of numerical data representing nodal values for a 450-triangle problem.

TABLE IV
Nodal Values of IRLSFEM Solution for Constant Convection Problem (450 Triangles)

Table with 15 columns and 15 rows of numerical data, mostly showing values of 2.0000000000000000.

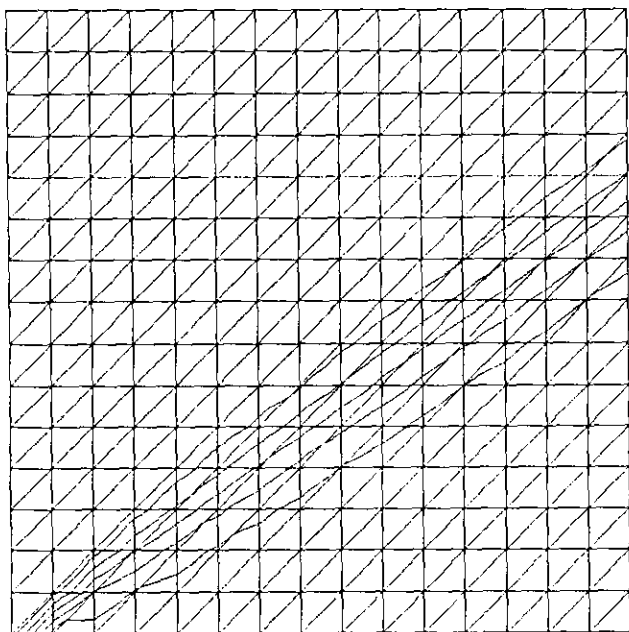


FIG. 2. Contours of LSFEM solution for constant convection problem (450 triangles).

double precision (8 bytes or 64 bits) real number in a computer can only represent a decimal number with 15 digits. This solution has no oscillation nor diffusion. The transition over the discontinuity is accurately located in the vicinity of the line $y = x \tan 35^\circ$ and is accomplished in just one element.

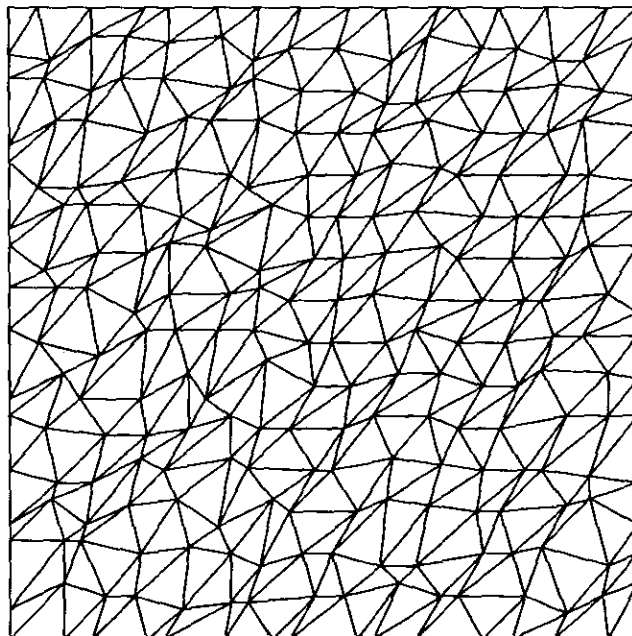


FIG. 4. Unstructured mesh with 450 triangles.

The least-squares solution for $n = 15$ (450 triangles) is given in Table III and depicted in Fig. 2. This solution is diffused and slightly oscillatory around the jump. Starting from this least-squares solution, the accurate discontinuous solution is obtained after four iterations (see Table IV and Fig. 3). This solution again has 15-digit accuracy. Because of the limitation of page size, we give the nodal values with only 8 digits in Table IV.

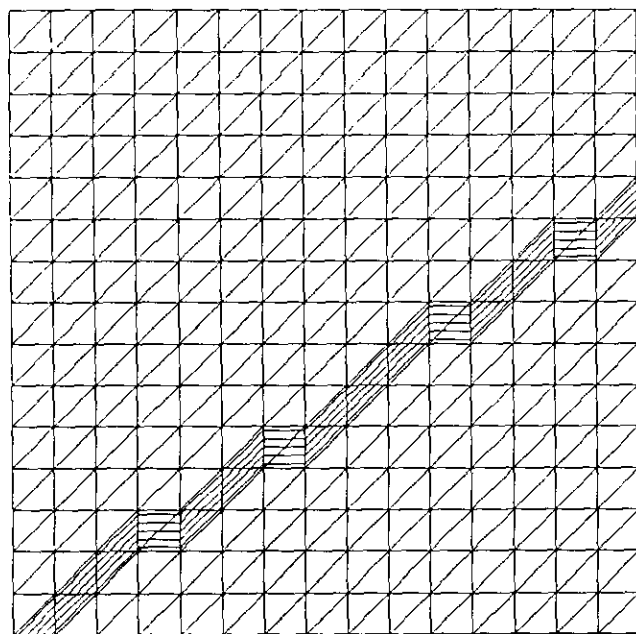


FIG. 3. Contours of IRLSFEM solution for constant convection problem (450 triangles).

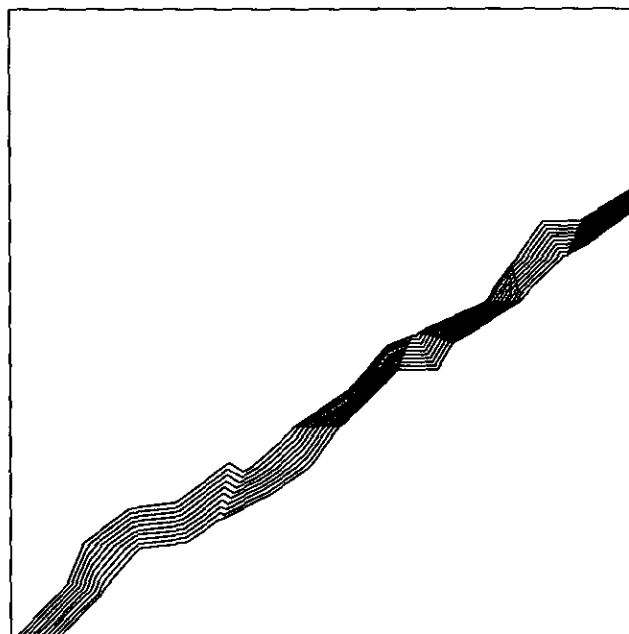


FIG. 5. Contours of IRLSFEM solution for constant convection problem on unstructured mesh.

TABLE V

Nodal Values of LSFEM Solution for Constant Convection Problem (5×5 Bilinear Elements)

2.00000000000000	2.00118821860769	2.01641923943302	2.03803285463046	2.02836549017373	1.94221664119897
2.00000000000000	2.01241750321268	2.02821882122379	2.00791191195100	1.90483206371286	1.72670330578520
2.00000000000000	2.01537346600061	1.99258728453429	1.86179903261828	1.62746979769939	1.37690702085225
2.00000000000000	1.98845751860405	1.82899110753578	1.51246025807770	1.25623579014445	1.07987201470759
2.00000000000000	1.83715624421649	1.34523292809047	1.12702336347438	1.02575100553365	0.98654845422136
2.00000000000000	1.00000000000000	1.00000000000000	1.00000000000000	1.00000000000000	1.00000000000000

TABLE VI

Nodal Values of IRLSFEM Solution for Constant Convection Problem (5×5 Bilinear Elements)

2.00000000000000	2.00000000000000	2.00000000000000	2.00000000000000	2.00000000000000	2.00000000000000
2.00000000000000	2.00000000000000	2.00000000000000	2.00000000000000	2.00000000000000	2.00000000000000
2.00000000000000	2.00000000000000	2.00000000000000	2.00000000000000	2.00000000000000	2.00000000000000
2.00000000000000	2.00000000000000	2.00000000000000	2.00000000000000	1.00000000000000	1.00000000000000
2.00000000000000	2.00000000000000	1.00000000000000	1.00000000000000	1.00000000000000	1.00000000000000
2.00000000000000	1.00000000000000	1.00000000000000	1.00000000000000	1.00000000000000	1.00000000000000

TABLE VII

Element Residuals of IRLSFEM Solution for Constant Convection Problem (5×5 Bilinear Elements)

0.18702E-14	0.17055E-15	0.49194E-14	0.48678E-14	0.22330E-14
0.92255E-15	0.49022E-14	0.35653E-14	0.11563E-15	0.18920E-14
0.18702E-14	0.18828E-14	0.59523E-16	0.20000E+01	0.20000E+01
0.85381E-15	0.20000E+01	0.20000E+01	0.20000E+01	0.14050E-14
0.20000E+01	0.20000E+01	0.49960E-15	0.64103E-16	0.49960E-15

TABLE VIII

Element Variations of IRLSFEM Solution for Constant Convection Problem (5×5 Bilinear Elements)

0.71054E-14	0.10658E-13	0.19540E-13	0.19540E-13	0.88818E-14
0.10658E-13	0.19540E-13	0.15987E-13	0.53291E-14	0.88818E-14
0.71054E-14	0.88818E-14	0.35527E-14	0.80000E+01	0.80000E+01
0.71054E-14	0.80000E+01	0.80000E+01	0.80000E+01	0.53291E-14
0.80000E+01	0.80000E+01	0.17764E-14	0.17764E-14	0.17764E-14

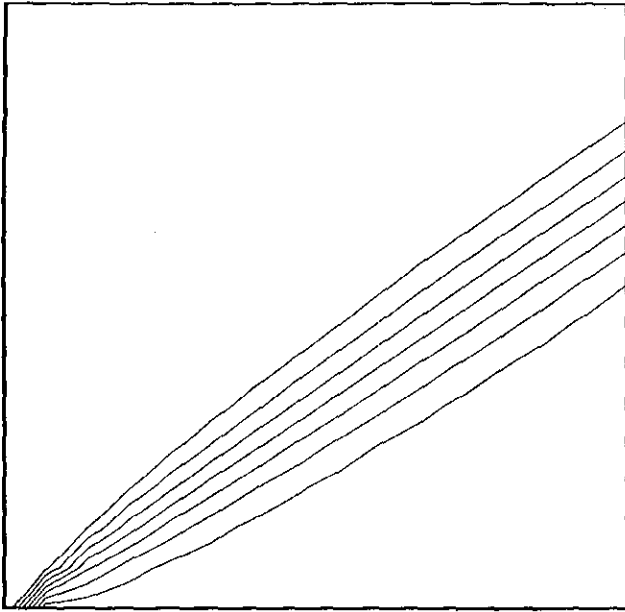


FIG. 6. Contours of LSFEM solution for constant convection problem (15×15 bilinear elements).

15×15 bilinear elements is presented in Fig. 6. This approximate solution is reasonably good, although the discontinuity is smeared out and slight oscillations occur. From this figure, we can hardly tell where the jump is located. However, after five iterations, a clean non-diffusive solution is reached (see Fig. 7). This solution again is 15-digit correct.

The least-squares solution on a mesh with 40×40 bilinear

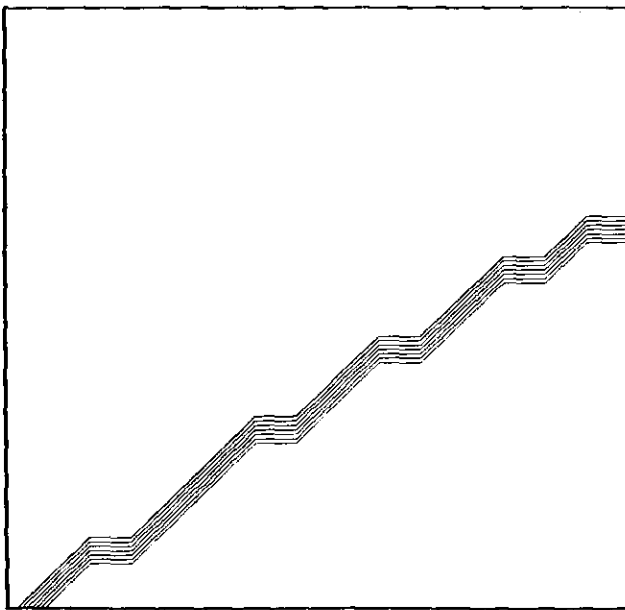


FIG. 7. Contours of IRLSFEM solution for constant convection problem (15×15 bilinear elements).

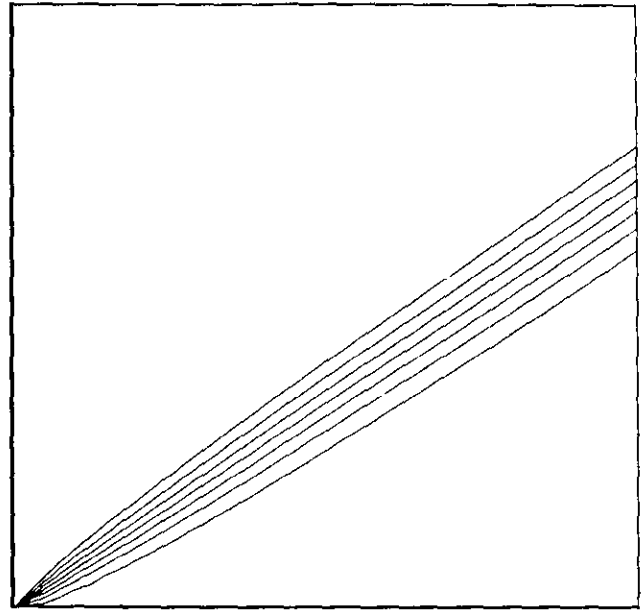


FIG. 8. Contours of LSFEM solution for constant convection problem (40×40 bilinear elements).

elements is illustrated in Fig. 8. Taking this initial least-squares solution, and after eight steps of processing, we obtained the highly accurate solution illustrated in Fig. 9.

We also did numerical tests for meshes with up to 100×100 elements, combined with various inflow angles and different boundary conditions. All of our IRLSFEM results are perfect. Because of the page limitation, we do not present these results here.

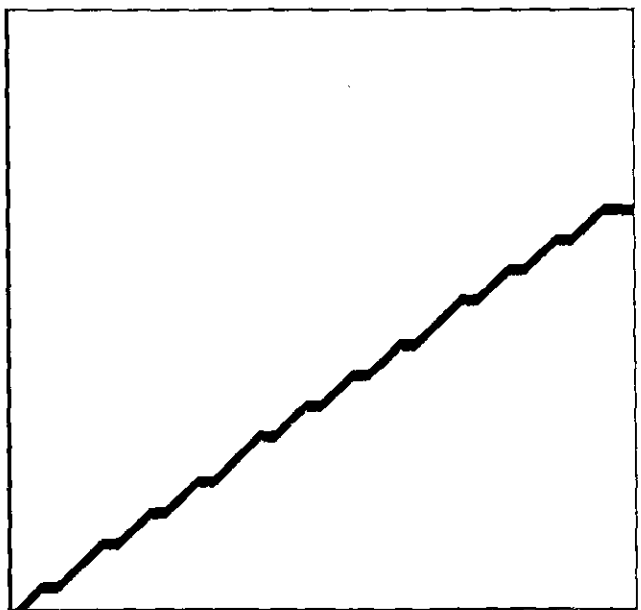


FIG. 9. Contours of IRLSFEM solution for constant convection problem (40×40 bilinear elements).

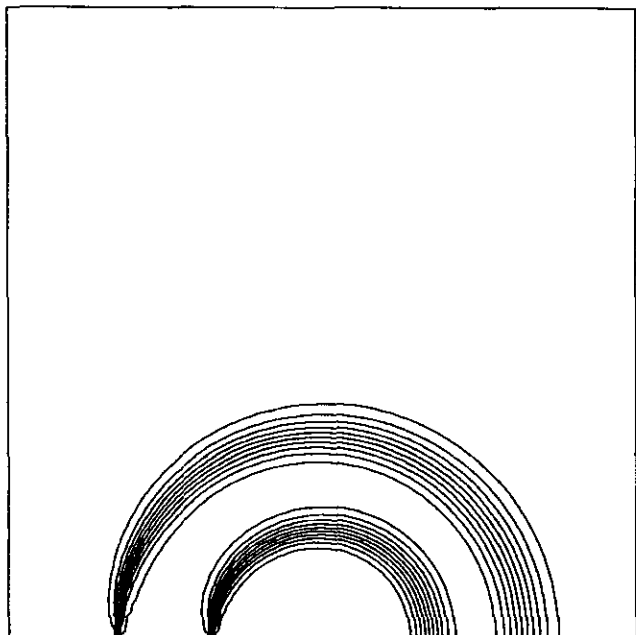


FIG. 10. Contours of LSFEM solution for circular convection problem (100 × 100 bilinear elements, 2 × 2 quadrature).

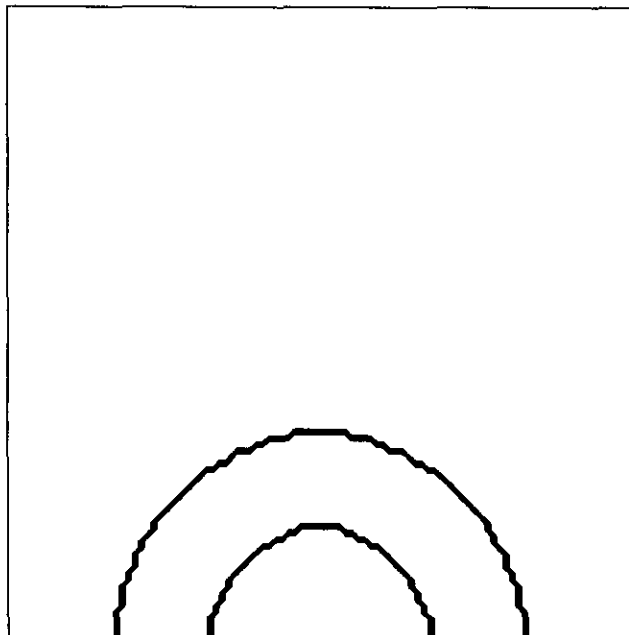


FIG. 12. Contours of IRLSFEM solution for circular convection problem (100 × 100 bilinear elements, 2 × 2 quadrature).

4.2. Spatially Varying Convection Field

Following Deconinck *et al.* [7], we considered the problem (1a) in the unit square with a circular convection field,

$$\beta_1 = y, \tag{27a}$$

$$\beta_2 = 0.5 - x, \tag{27b}$$

and boundary conditions,

$$u(0, y) = 0, \tag{27c}$$

$$u(x, 1) = 0, \quad x \geq 0.5 \tag{27d}$$

$$u(x, 0) = \begin{cases} 0 & x \leq 0.17; \\ 1 & 0.17 < x < 0.33; \\ 0 & 0.33 \geq x < 0.5. \end{cases} \tag{27e}$$

The mesh with 100 × 100 uniform bilinear elements was employed. The size of the elements in this mesh is 0.01

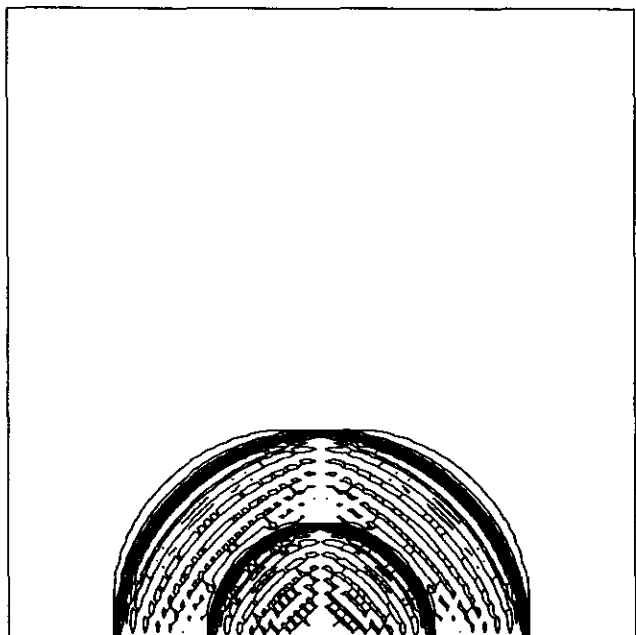


FIG. 11. Contours of LSFEM solution for circular convection problem (one-point quadrature).

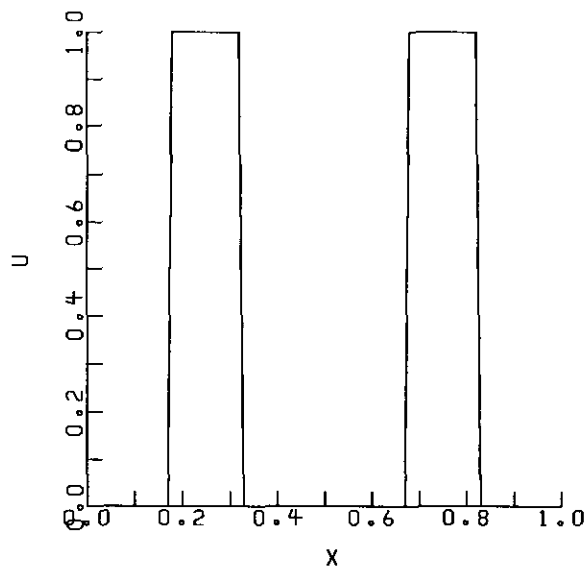


FIG. 13. Boundary distribution of u for circular convection problem.

(compare with 0.005 in the mesh used in [7]). Therefore, we modified the data a little bit in (27e) in order to make them consistent with the lower resolution of our mesh. In the least-squares finite element solution with 2×2 quadrature, as seen in Fig. 10, the input distribution (the left half of the lower boundary) is smeared quite a bit by the time it reaches the outflow (the right half of the lower boundary). Starting from this diffused solution, we still obtained a completely non-oscillatory and non-diffusive solution with correct jumps in just one element. However, the location of computed contact discontinuities on the outflow boundary has deviation of one or two grids from the exact solution. As indicated in Section 2, we may use the reduced integration (one point quadrature) to obtain a better initial solution, which is shown in Fig. 11. This solution looks terrible but contains the correct information about the location of discontinuities. Starting from this oscillatory solution, after eight reweighted iterations (with 2×2 quadrature) we obtained a perfect 14-digit correct discontinuous solution shown in Fig. 12. The distribution of u on the bottom boundary of the domain is shown in Fig. 13. No other currently available methods can produce such a sharp and highly accurate resolution.

5. DISCUSSION

In the early version of this paper [16], the author mistakenly considered that the correct discontinuous solution is the minimizer of the L_1 norm of the residual, and thus called the IRLSFEM an L_1 method. Immediately after the publication of [16], Lowrie and Roe [26] reproduced the results and found that the accurate discontinuous solution does not minimize the L_1 norm of the residual. They further proved that for problem (26) the L_1 minimizer and the accurate discontinuous solution cannot be the same. This is an important discovery. It also implies that our method is totally new from concept to implementation.

Obviously, the present version of the iteratively reweighted least-squares finite element method cannot be used for problems with smooth solutions such as problem (19). This method will take the elements with larger gradients as “shocked” elements and completely ignore them and will then yield a wrong discontinuous solution.

If the given data g in problem (1) are varying along the boundary Γ_- with large gradient and a small jump, we will have the same difficulty. This type of difficulty is related to the numerical measure of discontinuity which is different from the mathematical definition of discontinuity. That is, we must distinguish the discontinuity from large gradients. This type of difficulty can be overcome by using fine meshes and specifying an appropriate level of jump for a particular problem. Only if the variation indicator of an element defined in (24) exceeds this level, should the contribution of this element be eliminated.

The choice of the sixth power in (21) is based on our limited experience. In fact, the method is not sensitive to the number of the power. The third, fourth, fifth, seventh, ... powers, all of them, work well. The only difference is that the number of iterations to convergence is different. Here the essential issue is how to completely eliminate the equations in “shocked” elements, and it does not matter what means are employed. Of course, a theoretic investigation to find an optimal way is very welcome.

6. CONCLUSIONS

A new procedure, based on the iteration of least-squares finite element method for the solution of pure convection problems with contact discontinuities, is developed. The overdetermined algebraic system is inherently obtained by choosing an appropriate number of Gaussian points in the formation of element matrices. Through reweighting the contamination of “shocked” elements is eliminated.

This finite element method captures two-dimensional discontinuity in bands of elements that are only one element wide on both coarse and fine meshes. The solution of this method has no smearing nor oscillation and has superior accuracy. The method is simple and robust. The concept can also be applied to three-dimensional problems.

We believe that the methodology developed in this paper can be transferred into many other areas which deal with sharp fronts such as oil reservoir simulation, weather forecast, pollution control, and image enhancement. We have already tested this method for two-dimensional compressible flows with shocks. The preliminary results are encouraging.

ACKNOWLEDGMENTS

The author appreciates the discussion with Professor S. H. Chang, Dr. J. E. Lavery, Dr. T. L. Lin, and Dr. W. G. Xue in the preparation of this paper. The author also acknowledges the comments of Professor P. L. Roe and Mr. R. B. Lowrie.

REFERENCES

1. I. Barrodale and F. D. K. Roberts, *SIAM J. Numer. Anal.* **10**, 839 (1973).
2. G. F. Carey and B. N. Jiang, *Int. J. Numer. Methods Eng.* **26**, 81 (1988).
3. G. F. Carey and J. T. Oden, *Finite Elements: A Second Course, Vol. II* (Prentice-Hall, Englewood Cliffs, NJ, 1986), p. 199.
4. M. Chapman, *J. Comput. Phys.* **44**, 84 (1981).
5. I. Christie, D. F. Griffiths, A. R. Michell, and O. C. Ziekiewicz, *Int. J. Numer. Methods Eng.* **10**, 1389 (1976).
6. P. G. Ciarlet, *The Finite Element Method for Elliptic Problems* (North-Holland, Amsterdam, 1978).
7. H. Deconinck, K. G. Powell, P. L. Roe, and R. Struijs, AIAA-91-1532-CP.
8. J. E. Dendy, *SIAM J. Numer. Anal.* **11**, 637 (1974).

9. B. Engquist, P. Lötstedt, and B. Sjögreen, *Math. Comput.* **52**, 509 (1989).
10. C. A. J. Fletcher, *Computational Techniques for Fluid Dynamics 1, 2* (Springer-Verlag, Berlin, 1988).
11. A. Harten, *J. Comput. Phys.* **49**, 357 (1983).
12. A. Harten, *SIAM J. Numer. Anal.* **21**, 1 (1984).
13. C. Hirsch, *Numerical Computation of Internal and External Flows I* (Wiley, Chichester, 1988).
14. T. J. R. Hughes, *Int. J. Numer. Methods Fluids* **7**, 1261 (1987).
15. T. J. R. Hughes and A. Brooks, "A Multidimensional Upwind Scheme with No Crosswind Diffusion," in *Finite Element Methods for Convection Dominated Flows*, edited by T. J. Hughes, AMD, Vol. 34, (ASME, New York, 1979).
16. B. N. Jiang, NASA TM 103773, ICOMP-91-03.
17. B. N. Jiang and G. F. Carey, "Element-by-Element Preconditioned Conjugate Gradient Algorithm for Compressible Flow," in *Innovative Methods for Nonlinear Problems*, edited by W. K. Liu, T. Belytschko, and K. C. Park (Pineridge Press, Swansea, UK, 1984).
18. B. N. Jiang and G. F. Carey, *Int. J. Numer. Methods Fluids* **8**, 933 (1988).
19. B. N. Jiang and G. F. Carey, *Int. J. Numer. Methods Fluids* **10**, 557 (1990).
20. B. N. Jiang and L. A. Povinelli, *Comput. Methods Appl. Mech. Eng.* **81**, 13 (1990).
21. C. Johnson, *Numerical Solution of Partial Differential Equations by the Finite Element Method* (Cambridge Univ. Press, Cambridge, UK, 1987).
22. C. Johnson, U. Nävert, and J. Pitkäranta, *Comput. Methods Appl. Mech. Eng.* **45**, 285 (1984).
23. S. Koshizuka, C. B. Carrico, S. W. Lomperski, Y. Oka, and Y. Togo, *Comput. Mech.* **6**, 65 (1990).
24. J. E. Lavery, *J. Comput. Phys.* **79**, 436 (1988).
25. J. E. Lavery, *SIAM J. Numer. Anal.* **26**, 1081 (1989).
26. R. B. Lowrie and P. L. Roe, "On the Numerical Solution of Conservation Laws by Minimizing Residuals," Conservation Laws and Shock Capturing, University of South Carolina, Sep. 13, 1991.
27. K. W. Morton and A. K. Parrot, *J. Comput. Phys.* **36**, 249 (1980).
28. J. T. Oden and G. F. Carey, *Finite Elements: Mathematical Aspects, Vol. IV* (Prentice-Hall, Englewood Cliffs, NJ, 1983).
29. J. T. Oden and L. Demkowicz, "Advances in adaptive improvements: A survey of adaptive methods in computational fluid mechanics," in *State of the Art Surveys in Computational Mechanics*, edited by A. K. Noor and J. T. Oden (ASME, New York, 1988).
30. J. Peraire, M. Vahdati, K. Morgan and O. C. Zienkiewicz, *J. Comput. Phys.* **72**, 449 (1987).
31. O. Pironneau, *Finite Element Method for Fluids* (Wiley, London, 1989).
32. E. V. Vorozhtsov, *Comput. Fluids* **15**, 13 (1987).
33. E. V. Vorozhtsov, *Comput. Fluids* **18**, 35 (1990).
34. L. B. Wahlbin, *RAIRO* **8**, 109 (1974).