

SEMIDISCRETE LEAST SQUARES METHODS FOR LINEAR CONVECTION-DIFFUSION PROBLEM

TSU-FEN CHEN

Department of Mathematics, Box 19408, University of Texas at Arlington
 Arlington, Texas 76019, U.S.A.

(Received October 1991)

Abstract—In this paper, some approximate methods for solving linear convection-diffusion problems are presented. The methods consist in discretizing with respect to time and solving the resulting convection dominated elliptic problem for fixed time by least squares finite element methods. An analysis of least squares approximations is given, including optimal order estimates for piecewise polynomial approximation spaces. The model problem considered is the time-dependent convection dominated linear convection-diffusion equation. Numerical results for the Burgers' equation will also be presented.

1. INTRODUCTION

As a model for time-dependent convection-dominated linear convection-diffusion problem, consider the following problem: find the scalar function $u(x, t)$ such that

$$\begin{aligned} u_t + u_\beta - \varepsilon \Delta u &= f, & \text{in } \Omega \times I, \\ u(x, 0) &= u_0(x), & x \in \Omega, \\ u(x, t) &= 0, & x \in \Gamma, t \in I, \end{aligned} \tag{1.1}$$

where Ω is a bounded domain in \mathbb{R}^2 with boundary Γ , $u_t = \frac{\partial u}{\partial t}$, $u_\beta = \beta \cdot \nabla u$ with ∇ the gradient with respect to $x = (x_1, x_2) \in \mathbb{R}^2$, $\beta = (\beta_1, \beta_2)$ is a given smooth vector field, and $\varepsilon > 0$ is a small constant. Further, f and u_0 are given data, and $I = (0, T)$ is a given time interval. Note that for simplicity zero boundary data is considered here.

By replacing the time derivative in (1.1) by a backward-difference quotient, we define an approximate solution $u^k(x, t)$ for $t = nk = n\Delta t$, $n = 0, 1, 2, \dots$, by

$$\begin{aligned} \frac{u^k(x, t+k) - u^k(x, t)}{k} + u_\beta^k(x, t+k) - \varepsilon \Delta u^k(x, t+k) &= f, & \text{in } \Omega \times I, \\ u^k(x, 0) &= u_0(x), & x \in \Omega, \\ u_k(x, t) &= 0, & x \in \Gamma, t \in I. \end{aligned} \tag{1.2}$$

With $u^k(x, t) = v$, $u^k(x, t+k) = w$, we then have the following equation to solve for w , when v is known:

$$\begin{aligned} w + k w_\beta - k \varepsilon \Delta w &= v + k f(x, t+k), & \text{in } \Omega \times I, \\ w &= 0, & x \in \Gamma, t \in I. \end{aligned} \tag{1.3}$$

If we consider, instead of (1.2), the Crank-Nicolson formula

$$\frac{\tilde{u}^k(x, t+k) - \tilde{u}^k(x, t)}{k} + \frac{1}{2} \tilde{u}_\beta^k(x, t+k) - \varepsilon \frac{1}{2} \Delta \tilde{u}^k(x, t+k) + \frac{1}{2} \tilde{u}_\beta^k(x, t) - \varepsilon \frac{1}{2} \Delta \tilde{u}^k(x, t) = f,$$

This work is supported by the Research Enhancement Program of The University of Texas at Arlington.

Typeset by $\mathcal{A}\mathcal{M}\mathcal{S}\text{-}\mathcal{T}\mathcal{E}\mathcal{X}$

the problem (1.3) changes into

$$\begin{aligned} w + \frac{1}{2}k w_\beta - \frac{1}{2}k \varepsilon \Delta w &= v - \frac{1}{2}k v_\beta + \frac{1}{2}k \varepsilon \Delta v + kf(x, t + \frac{1}{2}k), & \text{in } \Omega \times I, \\ w &= 0, & x \in \Gamma, t \in I. \end{aligned} \quad (1.4)$$

Note that when $\varepsilon > 0$ and k is bounded away from zero, the problems (1.3) and (1.4) admit a unique solution u^k and \tilde{u}^k , respectively, and it can be shown that, for sufficiently smooth initial values v ,

$$\begin{aligned} \sup_{0 \leq t \leq T} \|u^k(\cdot, t) - u(\cdot, t)\| &= O(k) \text{ as } k \rightarrow 0, \text{ and} \\ \sup_{0 \leq t \leq T} \|\tilde{u}^k(\cdot, t) - u(\cdot, t)\| &= O(k^2) \text{ as } k \rightarrow 0, \end{aligned} \quad (1.5)$$

where $\|\cdot\|$ denotes the norm in $L_2(\Omega)$, $\|v\| = \left(\int_{\Omega} |v(x)|^2 dx \right)^{1/2}$.

When the Galerkin method is applied to (1.3) or (1.4), it may produce an oscillatory solution if $\varepsilon < h$, where h is the spatial mesh size and the exact solution is not smooth. The difficulty is that the method only works well for equations which are diffusion dominated, in the sense that $\varepsilon > h$. In particular, the error made in using the Galerkin method (with linear test functions) for this problem is $O(h^2/\varepsilon)$, which indicates that the method is ineffective for very small ε . One possible solution is to replace ε with $\tilde{\varepsilon}$, which will keep the equation diffusion-dominated. This corresponds to the backward spatial differencing in the finite difference calculations. In the context of finite element methods, test functions were selectively modified in the convective term ('upwinding') to simulate the directional properties of the reduced hyperbolic problem. Further generalization to this approach were developed and led to the Petrov-Galerkin methods which uses different finite dimensional subspaces for the trial and test functions [1-5]. However, there are several detractors to these methods. The most significant is the lack of a systematic procedure to extend the methods to more general linear and nonlinear problems. The need to optimize certain parameters to control oscillations and dissipation is also a limitation.

Semidiscrete least squares methods to the heat equation were considered by Bramble and Thomée [6]. They approximated the solution of the heat equation by minimizing the L_2 norm of the residual functionals obtained from (1.3) and (1.4). Restrictions on the relation between the spatial mesh size h and $k (= \Delta t)$ were needed in their analysis. In addition, with the presence of the diffusion term $-\varepsilon \Delta w$, linear test functions can not be used in the approximations. In the following, we describe the least squares method which allows the use of the linear test functions. Throughout the paper, only the Crank-Nicolson scheme (1.4) will be considered. The analysis of the purely implicit scheme (1.3) follows directly from that of (1.4). To approximate (1.4), for simplicity, we assume that β is a constant vector and let

$$u = w\beta - \varepsilon \nabla w.$$

Thus,

$$\begin{aligned} u - w\beta + \varepsilon \nabla w &= 0, \\ w + \frac{1}{2}k \operatorname{div} u &= v - \frac{1}{2}k \operatorname{div} z, \quad \text{where } z = v\beta - \varepsilon \nabla v. \end{aligned} \quad (1.6)$$

Our approximation is then obtained by minimizing

$$\int_{\Omega} \left\{ |u - w\beta + \varepsilon \nabla w|^2 + \left| w + \frac{1}{2}k \operatorname{div} u - \left(v - \frac{1}{2}k \operatorname{div} z \right) \right|^2 \right\} \quad (1.7)$$

over all u and w with appropriate boundary conditions. This is referred as the least squares (LS) method. We also consider minimizing a variation of (1.7) (referred as the weighted least squares (WDLS) method)

$$\int_{\Omega} \left\{ \chi |u - w\beta + \varepsilon \nabla w|^2 + \left| w + \frac{1}{2}k \operatorname{div} u - \left(v - \frac{1}{2}k \operatorname{div} z \right) \right|^2 \right\} \quad (1.8)$$

where χ is a positive weighting function.

Variational principles of the least squares types have a number of valuable computational properties. For example, the algebraic system generated is always Hermitian semidefinite. In addition, such schemes, if properly formulated, are insensitive to the type of the partial differential equation, i.e., the computational algorithm is the same in elliptic and hyperbolic regions. Successful applications to the transonic flow problems based on the formulation were presented in [7] and [8]. In addition, an analysis of the least squares approximation to elliptic boundary value problems were presented in [7]. Based on the theoretical framework, the stability and error estimates of the method can be established for (1.7) and will be discussed.

Following this introduction, the remaining is divided into four sections. Section 2 introduces the variational formulation based on (1.7) and the necessary assumptions which includes a special grid decomposition property introduced in [9]. In Section 3, we will provide an error analysis for the least squares (LS) method which gives the optimal error estimates for both u and w , if a special regularity property is satisfied and $k = Ch$. In the case when ε is close to zero, $u \approx w\beta$ and the estimate in u does not reflect how well ∇w is approximated. Thus the weighted least squares method (WDLS) with appropriate weighting function χ will be introduced and analyzed in Section 4. With the proper choice of the weighting function χ , unlike the (LS) method, this new formulation will give optimal order error in both w and u without using the grid decomposition property. Finally, in Section 5, numerical results for the Burgers' equation will be presented and compared to the upwinded method in [1]. In order to handle the boundary layer effect when ε is close to zero, a variation of (LS) and (WDLS) will also be introduced in Section 5. It is demonstrated that this new approach acts naturally in a manner similar to upwinding and requires no "free" parameters.

2. FORMULATION OF PROBLEM AND ASSUMPTIONS

For clarity, instead of (1.5), we consider least squares approximation to the following: given a function g and step size $k = \Delta t$, we seek a suitably smooth function ϕ satisfying

$$\phi + \frac{1}{2} k \phi_\beta - \frac{1}{2} k \varepsilon \nabla \phi = g, \quad \text{in } \Omega, \quad (2.1)$$

$$\phi = 0, \quad \text{on } \Gamma, \quad (2.2)$$

or what is the same

$$u - \phi \beta + \varepsilon \nabla \phi = 0, \quad \text{in } \Omega, \quad (2.3)$$

$$\phi + \frac{1}{2} k \operatorname{div} u = g, \quad \text{in } \Omega, \quad (2.4)$$

$$\phi = 0, \quad \text{on } \Gamma. \quad (2.5)$$

Note that, for simplicity, we assume that β is a constant vector. When $\beta = (\beta_1, \beta_2)$ is not a constant vector, we still have $u - \phi \beta + \varepsilon \nabla \phi = 0$ and (2.4) becomes $\phi + \frac{1}{2} k \operatorname{div} u - \frac{1}{2} k \phi \operatorname{div} \beta = g$. To be precise, we assume $g \in L_2(\Omega)$ and we seek solution ϕ, u to (2.3)–(2.5) in

$$S_1 = \{\psi | \psi \in H^1(\Omega), \psi = 0 \text{ on } \Gamma\} \quad \text{on} \quad V_0 = \{v | v \in H^1(\Omega)\}.$$

Let $\|\cdot\|, |\cdot|$ be norms on V_0, S_1 with $(\cdot, \cdot), \langle \cdot, \cdot \rangle$ being inner products. Here $\|\cdot\|$ and $|\cdot|$ are L_2 norms.

To approximate, we introduce finite dimensional subspaces

$$S_h \subseteq S_1, \quad V_h \subseteq V_0.$$

We determine $u_h \in V_h, \phi_h \in S_h$ by minimizing

$$\|v_h - \psi_h \beta + \varepsilon \nabla \psi_h\|^2 + \left| \psi_h + \frac{1}{2} k \operatorname{div} v_h - g \right|^2 \quad (2.6)$$

over v_h in V_h and ψ_h in S_h . Taking the first variation gives

$$\begin{aligned} & (u_h - \phi_h \beta + \varepsilon \nabla \phi_h, v^h - \psi^h \beta + \varepsilon \nabla \psi^h) + \left\langle \phi_h + \frac{1}{2} k \operatorname{div} u_h, \psi^h + \frac{1}{2} k \operatorname{div} v^h \right\rangle \\ & = \left\langle g, \psi^h + \frac{1}{2} k \operatorname{div} v^h \right\rangle, \end{aligned} \quad (2.7)$$

a relation which holds for all $v^h \in V_h$ and $\psi^h \in S_h$. A useful fact is that (2.7) remains valid when ϕ_h is replaced by ϕ and u_h is replaced by u , where $\{\phi, u\}$, $\phi \in S_1$, $u \in V_0$ is the solution of (LS), i.e.,

$$\begin{aligned} & (u - \phi \beta + \varepsilon \nabla \phi, v^h - \psi^h \beta + \varepsilon \nabla \psi^h) + \left\langle \phi + \frac{1}{2} k \operatorname{div} u, \psi^h + \frac{1}{2} k \operatorname{div} v^h \right\rangle \\ & = \left\langle g, \psi^h + \frac{1}{2} k \operatorname{div} v^h \right\rangle. \end{aligned} \quad (2.8)$$

We shall assume throughout the standard approximation properties [10] for the finite dimensional spaces V_h, S_h in terms of the Sobolev norms $\|\cdot\|_r$ on $H^r(\Omega)$. In particular, we shall need the following assumptions.

2.1. Approximation Property

For any u and ϕ in $H^r(\Omega)$, there exist interpolants $\tilde{u}^h \in V_h$ and $\tilde{\phi}^h \in S_h$ such that

$$\|u - \tilde{u}^h\|_l \leq C_A h^{r-l} \|u\|_r, \quad \text{and} \quad \|\phi - \tilde{\phi}^h\|_l \leq C_A h^{r-l} \|\phi\|_r \quad (2.9)$$

for $l = 0$ and $l = 1$ where C_A is a constant independent of h, u and ϕ .

Another assumption used in the theory is the Grid Decomposition Property introduced in [10]. A precise statement of this property is as follows.

2.2. Grid Decomposition Property (GDP)

For each $v_h \in V_h$ there exists w_h and z_h in V_h for which

$$v_h = w_h + z_h$$

with $\operatorname{div} z_h = 0$ and

$$(z_h, w_h) = \int_{\Omega} z_h \cdot w_h = 0, \quad \|w_h\| = \|w_h\|_0 \leq C_G \|\operatorname{div} v_h\|_{-1} \quad (2.10)$$

for some positive constant C_G independent of h and v_h . Note that if triangular linear elements are used, this property will not hold for the directional grids but it is valid for the criss-cross grids in Figure 1 [9].

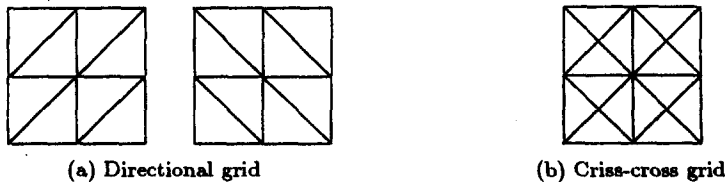


Figure 1.

3. ERROR ANALYSIS FOR THE LEAST SQUARES (LS) METHOD

The error analysis of least-squares methods starts from the observation that the solution $\{\phi_h, u_h\}$ of the discrete problem is a best approximation to $\{\phi, u\}$ in a suitable norm. This norm arises naturally from the bilinear form on $S_1 \times V_0$

$$a((\phi, u), (\psi, v)) = (u - \phi\beta + \varepsilon \nabla\phi, v - \psi\beta + \varepsilon \nabla\psi) + \left\langle \phi + \frac{1}{2}k \operatorname{div} u, \psi + \frac{1}{2}k \operatorname{div} v \right\rangle \quad (3.1)$$

and is given by

$$|||(\phi, u)||| = a((\phi, u), (\phi, u))^{1/2}.$$

Let $\eta = \phi - \phi_h$, $e = u - u_h$. We have (2.7) and (2.8) that imply the error $\{\eta, e\}$ is orthogonal to $S_h \times V_h$ in the form $a(\cdot, \cdot)$; i.e.,

$$a((\eta, e), (\psi^h, v^h)) = 0 \quad \forall (\psi^h, v^h) \in S_h \times V_h. \quad (3.2)$$

Observe that, with $|\beta|$, being an upper bound for the vector β ,

$$|||(\psi, v)||| \leq \left\{ \|v\|_0 + \frac{1}{2}k\|v\|_1 + (1 + |\beta|)\|\psi\|_0 + \varepsilon\|\psi\|_1 \right\}. \quad (3.3)$$

This follows immediately from (3.1) and the fact that $a(\cdot, \cdot)$ is a bounded form on $H^1(\Omega) \times H^1(\Omega)$.

To simplify the notations, we introduce the following. Let

$$e_h(v) = \inf_{v^h \in V_h} \|v - v^h\|$$

denote the error in the best approximation. Note that from the approximation property (2.9)

$$e_h(v) \leq C_A h^r \|v\|_r.$$

In our theory, we have two rather complicated error terms σ and γ_h :

$$\sigma_h = \sup_{\xi} \left\{ \sigma_1(\xi) \mid \|\xi\|_2 \leq 1 \right\}, \quad \gamma_h = \sup_{\xi} \left\{ \gamma_1(\xi, \xi\beta - \varepsilon \nabla\xi) \mid \|\xi\|_3 \leq 1 \right\}.$$

where

$$\sigma_1(\xi) = \inf_{\xi \in S_h} \left\{ \left\| \xi_h + \frac{2}{k}\xi \right\|_0 + \|\nabla\xi + \xi_h\beta - \varepsilon \nabla\xi_h\|_0 \right\}$$

and

$$\gamma_1(\xi, y) = \inf_{(\xi_h, y_h) \in S_h \times V_h} \left\{ \|y - y_h\|_0 + \frac{1}{2}k\|y - y_h\|_1 + (1 + |\beta|)\|\xi - \xi_h\|_0 + \varepsilon\|\xi - \xi_h\|_1 \right\}.$$

If S_h and V_h consist of piecewise linear elements, then

$$\sigma_h \leq Ch, \quad \gamma_h \leq Ch,$$

follows directly from approximation property (2.9).

The following is an immediate consequence of the orthogonality (3.2) and implies that $\{\phi_h, u_h\}$ is a best approximation to $\{\phi, u\}$ in $||| \cdot |||$.

LEMMA 3.1.

$$|||(\eta, e)||| \leq |||(\phi - \psi^h, u - v^h)|||, \quad \text{for all } (\psi^h, v^h) \in S_h \times V_h. \quad (3.4)$$

Lemma 3.1 implies that

$$|||(\eta, e)||| \leq C_A \left\{ \left(h^r + \frac{1}{2}k h^{r-1} \right) \|u\|_r + \left((1 + |\beta|)h^r + \varepsilon h^{r-1} \right) \|\phi\|_r \right\}, \quad (3.5)$$

which is $O(h^{r-1})$. Note that this error estimate is not very useful by itself since the reverse of (3.3) is not valid, i.e., $||| \cdot |||$ is majorized by the norm $\|\cdot\|_1$ on $H^1(\Omega)$ but is not equivalent to it. However, we can use it and the solvability of the dual problem of (2.1)–(2.2) to establish the optimal error estimate for $|\eta|$.

THEOREM 3.2.

$$|\eta| \leq C_R \sigma_h \| |(\eta, e) | \| . \quad (3.6)$$

PROOF. We solve for $\zeta \in S_1$ satisfying

$$\zeta - \frac{1}{2} k \zeta \beta - \frac{1}{2} k \varepsilon \Delta \zeta = \frac{1}{2} k \eta \quad \text{in } \Omega, \quad \zeta = 0 \quad \text{on } \Gamma. \quad (3.7)$$

Note that from the approximation theory [10], there is a positive constant C_R such that

$$\|\zeta\|_2 \leq C_R \|\eta\|_0.$$

We then solve for $\xi \in S_1$ such that

$$\frac{1}{2} k \operatorname{div} (\xi \beta - \varepsilon \nabla \xi) + \xi = -\frac{1}{2} k \Delta \zeta + \frac{2}{k} \zeta \quad \text{in } \Omega, \quad \xi = 0 \quad \text{on } \Gamma.$$

Note that

$$a((\eta, e), (\xi + \xi_h, \nabla \zeta + \xi \beta - \varepsilon \nabla \xi)) = \frac{2}{k} \left\langle \eta + \frac{1}{2} k \operatorname{div} e, \zeta \right\rangle + (\varepsilon - \eta \beta + \varepsilon \nabla \eta, \nabla \zeta).$$

Using the Green's identity $(\nabla \eta, \nabla \zeta) = -(\eta, \nabla \zeta)$ and $(e, \nabla \zeta) = -(\operatorname{div} e, \zeta)$, we have

$$(e - \eta \beta + \varepsilon \nabla \eta, \nabla \zeta) = -\frac{2}{k} \left\langle \eta + \frac{1}{2} k \operatorname{div} e, \zeta \right\rangle + (\eta, \eta).$$

Thus,

$$a((\eta, e), (\xi + \xi_h, \nabla \zeta + \xi \beta - \varepsilon \nabla \xi)) = |\eta|^2,$$

and then

$$|\eta|^2 \leq \| |(\eta, e) | \| \| |(\xi + \xi_h, \nabla \zeta + \xi \beta - \varepsilon \nabla \xi) | \|.$$

Therefore, we obtain (3.6) since

$$\| |(\xi + \xi_h, \nabla \zeta + \xi \beta - \varepsilon \nabla \xi) | \| = \left\{ \|\xi_h + \frac{2}{k} \zeta\|_0^2 + \|\nabla \zeta + \xi_h \beta - \varepsilon \nabla \xi_h\|_0^2 \right\}^{1/2} \leq C_R \sigma_h |\eta|.$$

To obtain error estimates for $\|e\|$, we shall need to exploit the solvability of the boundary value problem (2.1)–(2.2). More precisely, we shall need an *a priori* inequality of the form

$$\|\psi\|_{2+l} < C_E \left\| \psi + \frac{1}{2} k \psi \beta - \frac{1}{2} k \varepsilon \Delta \psi \right\|_l, \quad l = 0, 1, \quad (3.8)$$

to hold for all $\psi \in H^{2+l}(\Omega)$ satisfying the boundary condition $\psi = 0$ on Γ . This will be the case for a fixed positive number C_E provided Ω and β are sufficiently smooth [11]. Note that C_E will vary with k . This regularity property will enable us to establish optimal error estimates for

$$\left\| \eta + \frac{1}{2} k \operatorname{div} e \right\|_{-1}.$$

In addition, it is essential to use the Grid Decomposition Property discussed in Section 2 to establish optimal error in $\|e\|$.

LEMMA 3.3. Let $C_D = k \varepsilon C_E$, where C_E is the constants in (3.8), then

$$\left\| \eta + \frac{1}{2} k \operatorname{div} e \right\|_{-1} \leq C_D \gamma_h \| |(\eta, e) | \| . \quad (3.9)$$

PROOF. We recall that

$$\left\| \eta + \frac{1}{2} k \operatorname{div} e \right\|_{-1} \leq \sup_{\theta} \left\{ \left| \left\langle \eta + \frac{1}{2} k \operatorname{div} e, \theta \right\rangle \right| \mid \|\theta\|_1 \leq 1 \right\}.$$

To prove (3.9), let $\theta \in S_1$ be given with $\|\theta\|_1 \leq 1$ and solve for $\xi \in S_1$ satisfying

$$\xi + \frac{1}{2} k \xi \beta - \frac{1}{2} k \varepsilon \Delta \xi = k \varepsilon \theta \quad \text{in } \Omega, \quad \xi = 0 \quad \text{on } \Gamma, \quad (3.10)$$

or what is the same

$$\begin{aligned} y - \xi \beta + \varepsilon \nabla \xi &= 0, & \text{in } \Omega, \\ \xi + \frac{1}{2} k \operatorname{div} y &= k \varepsilon \theta, & \text{in } \Omega, \\ \xi &= 0, & \text{on } \Gamma. \end{aligned}$$

Regularity gives

$$\|\xi\|_3 \leq C_E k \varepsilon \|\theta\|_1 \leq C_D \|\theta\|_1.$$

Note that

$$a((\eta, e), (\xi - \xi_h, y - y_h)) = \left\langle \eta + \frac{1}{2} k \operatorname{div} e, \theta \right\rangle, \quad \text{all } (\xi_h, y_h) \in S_h \times V_h.$$

Therefore,

$$\left| \left\langle \eta + \frac{1}{2} k \operatorname{div} e, \theta \right\rangle \right| \leq \|(\eta, e)\| \|(\xi - \xi_h, y - y_h)\|.$$

Taking the infimum over $(\xi_h, y_h) \in S_h \times V_h$ and then taking the supremum over $\theta \in S_1$ with $\|\theta\|_1 \leq 1$ gives (3.9).

Finally, we complete the error analysis, making the use of the Grid Decomposition Property.

LEMMA 3.4. *Let the Grid Decomposition Property (GDP) hold and let \tilde{u}_h be a best approximation in the sense that*

$$\|u - \tilde{u}_h\| = e_h(u).$$

Then

$$\|u_h - \tilde{u}_h\| \leq 2C_G \left\| \frac{1}{2} k \operatorname{div} (u_h - \tilde{u}_h) \right\|_{-1} + e_h(u) + |\beta| |\eta|.$$

PROOF. We use (GDP) to write

$$\frac{1}{2} k (u_h - \tilde{u}_h) = w_h + z_h \quad \text{with } \operatorname{div} z_h = 0 \text{ and}$$

$$(z_h, w_h) = \int_{\Omega} z_h \cdot w_h = 0, \quad \|w_h\| = \|w_h\|_0 \leq C_G \left\| \frac{1}{2} k \operatorname{div} (u_h - \tilde{u}_h) \right\|_{-1}.$$

Note that from the orthogonality (3.2), we have

$$a((\phi, u - \tilde{u}_h), (\psi^h, v^h)) = a((\phi_h, u_h - \tilde{u}_h), (\psi^h, v^h)), \quad \text{all } (\psi^h, v^h) \in S_h \times V_h.$$

Choosing $\psi^h \in S_h$ such that,

$$\psi^h + \frac{1}{2} k \operatorname{div} (\psi^h \beta - \varepsilon \nabla \psi^h) = \operatorname{div} z_h \quad \text{and let } v^h = z_h + \psi^h \beta - \varepsilon \nabla \psi^h,$$

we have

$$(u - \tilde{u}_h - (\phi - \phi_h) \beta, z_h) = (u_h - \tilde{u}_h, z_h).$$

Then

$$\|z_h\| \leq e_h(u) + |\beta| |\eta| + \|w_h\|.$$

THEOREM 3.5. *Let the (GDP) hold. Then*

$$\|u - u_h\| \leq (2 + C_G)e_h(u) + C\|(\eta, e)\|, \quad (3.11)$$

where

$$C = 2C_D C_G \gamma_h + (2C_G + |\beta|) C_R \sigma_h.$$

PROOF. Note that

$$\left\| \frac{1}{2} k \operatorname{div} (u_h - \tilde{u}_h) \right\|_{-1} \leq \left\| \frac{1}{2} k \operatorname{div} (u - u_h) \right\|_{-1} + \left\| \frac{1}{2} k \operatorname{div} (u - \tilde{u}_h) \right\|_{-1}.$$

But from the definition of $\|\cdot\|_{-1}$,

$$\left\| \frac{1}{2} k \operatorname{div} (u - \tilde{u}_h) \right\|_{-1} \leq \frac{1}{2} k \|u - \tilde{u}_h\|.$$

Also, using Lemma 2, we have

$$\left\| \frac{1}{2} k \operatorname{div} (u - u_h) \right\|_{-1} \leq \left\| \eta + \frac{1}{2} k \operatorname{div} e \right\|_{-1} + \|\eta\|_{-1} \leq C_D \gamma_h \|(\eta, e)\| + |\eta|.$$

Combining these results with Lemma 3.4, we obtain (3.11).

Note that the estimate in $\varepsilon \|\nabla \phi\|$ is optimal which follows from the optimal estimate in u and ϕ . However, this does not reflect the accuracy of the approximation in $\|\nabla \phi\|$ when ε is close to zero. Even in the case when $\varepsilon > h$, the estimate (3.6) in η lies heavily on the solvability of the dual problem (3.7). In addition, the nonvanishing term $\varepsilon \nabla \phi$ plays an essential part of the analysis. As for the estimate (3.11) in u , it depends on the extra regularity property (3.8) and the validation of the Grid Decomposition Property (2.10). Based on these considerations, in the following section we introduce a weighted least squares formulation which gives the optimal estimate in ϕ without using the solvability of the dual problem (3.7). Moreover, a better estimate of $\|\nabla \phi\|$ will be obtained with this new formulation when ε is close to zero.

4. ERROR ANALYSIS FOR THE WEIGHTED LEAST SQUARE (WDLS) METHOD

In this section, we consider the following weighted least squares (WDLS) method. Slightly different from the (LS) in Section 3, we determine $u_h \in V_h$, $\phi_h \in S_h$ by minimizing

$$\frac{kh}{2\varepsilon\delta} \|v_h - \psi_h \beta + \varepsilon \nabla \psi_h\|^2 + \left| \psi_h + \frac{1}{2} k \operatorname{div} v_h - g \right|^2 \quad (4.1)$$

over v_h in V_h and ψ_h in S_h , where δ is a constant such that $\delta \leq h$. The first variation of (4.1) gives

$$\begin{aligned} & \frac{kh}{2\varepsilon\delta} (u_h - \phi_h \beta + \varepsilon \nabla \phi_h, v^h - \psi^h \beta + \varepsilon \nabla \psi^h) + \left\langle \phi_h + \frac{1}{2} k \operatorname{div} u_h, \psi^h + \frac{1}{2} k \operatorname{div} v^h \right\rangle \\ & = \left\langle g, \psi^h + \frac{1}{2} k \operatorname{div} v^h \right\rangle, \end{aligned} \quad (4.2)$$

a relation which holds for all $v^h \in V_h$ and $\psi^h \in S_h$. Note that (4.2) remains valid when ϕ_h is replaced by ϕ and u_h is replaced by u , where $\{\phi, u\}$, $\phi \in S_1$, $u \in V_0$ is the solution of (WDLS), i.e.,

$$\begin{aligned} & \frac{kh}{2\varepsilon\delta} (u - \phi \beta + \varepsilon \nabla \phi, v^h - \psi^h \beta + \varepsilon \nabla \psi^h) + \left\langle \phi + \frac{1}{2} k \operatorname{div} u, \psi^h + \frac{1}{2} k \operatorname{div} v^h \right\rangle \\ & = \left\langle g, \psi^h + \frac{1}{2} k \operatorname{div} v^h \right\rangle. \end{aligned}$$

For simplicity, we define a bilinear form on $S_1 \times V_0$

$$b((\phi, u), (\psi, v)) = \frac{kh}{2\varepsilon\delta} (u - \phi\beta + \varepsilon\nabla\phi, v - \psi\beta + \varepsilon\nabla\psi) + \left\langle \phi + \frac{1}{2}k \operatorname{div} u, \psi + \frac{1}{2}k \operatorname{div} v \right\rangle.$$

Further we let

$$|||(\psi, v)|||_b = b((\psi, v), (\psi, v))^{1/2}.$$

Observe that, with β being an upper bound for the vector β ,

$$|||(\psi, v)|||_b \leq \left\{ \left(\frac{kh}{2\varepsilon\delta} \right)^{1/2} \|v\|_0 + \frac{1}{2}k \|v\|_1 + \left(1 + |\beta| \left(\frac{kh}{2\varepsilon\delta} \right)^{1/2} \right) \|\psi\|_0 + \left(\frac{kh\varepsilon}{2\delta} \right)^{1/2} \|\psi\|_0 \right\}. \quad (4.3)$$

Let $\eta = \phi - \phi_h$, $e = u - u_h$. We then have the ‘‘orthogonality’’ which is similar to (3.2)

$$b((\eta, e), (\psi^h, v^h)) = 0 \quad \text{all } (\psi^h, v^h) \in S_h \times V_h. \quad (4.4)$$

We proceed with the error estimates by showing that $\{\phi_h, u_h\}$ of the discrete problem is a best approximation to $\{\phi, u\}$ in $|||(\cdot, \cdot)|||_b$ norm.

LEMMA 4.1.

$$|||(\eta, e)|||_b \leq |||(\phi - \psi^h, u - v^h)|||_b \quad \text{for all } (\psi^h, v^h) \in S_h \times V_h. \quad (4.5)$$

PROOF. This follows immediately from (4.4).

Lemma 4.1 implies that if $k = O(h)$ and linear elements are used in both S_h and V_h ,

$$|||(\eta, e)|||_b \leq Ch^2. \quad (4.6)$$

Note that (4.6) is of optimal order. Thus, different than the similar orthogonality result in Section 2, the estimate (4.5) helps us to establish the following results. We now use (4.5) to establish error estimates for $|\eta|$, $|||\nabla\eta|||$ and $|\operatorname{div} e|$.

THEOREM 4.2. If $0 < C_1 \leq k \leq C_2 h$ and $\delta = C_2/C_1$, then

$$|\eta| \leq |||(\eta, e)|||_b, \quad (4.7)$$

$$|||\nabla\eta||| \leq \sqrt{2} \left(\frac{\delta}{C_2\varepsilon} \right)^{1/2} |||(\eta, e)|||_b, \quad \text{and} \quad (4.8)$$

$$|\operatorname{div} e| \leq \frac{2\delta}{C_2} |||(\eta, e)|||_b. \quad (4.9)$$

PROOF. For any $(\psi, v) \in S_1 \times V_0$, since $\psi = 0$ on Γ , it follows that

$$\begin{aligned} |||(\psi, v)|||_b^2 &= b((\psi, v), (\psi, v)) \geq \frac{k}{2\varepsilon} \|v - \psi\beta + \varepsilon\nabla\psi\|^2 + \left| \psi + \frac{1}{2}k \operatorname{div} v \right|^2 \\ &= \frac{k}{2\varepsilon} \|v - \psi\beta\|^2 + \frac{1}{2}k\varepsilon \|\nabla\psi\|^2 + |\psi|^2 + \frac{1}{4}k^2 |\operatorname{div} v|^2. \end{aligned} \quad (4.10)$$

Note that $(\eta, e) \in S_1 \times V_0$, thus (4.10) is valid for (η, e) . Therefore the theorem follows directly from (4.10).

THEOREM 4.3. If $0 < C_1 \leq k \leq C_2 h$ and $\delta = C_2/C_1$, then

$$\|e\| \leq \left\{ \left(\frac{2\varepsilon\delta}{C_2} \right)^{1/2} + |\beta| \right\} |||(\eta, e)|||_b. \quad (4.11)$$

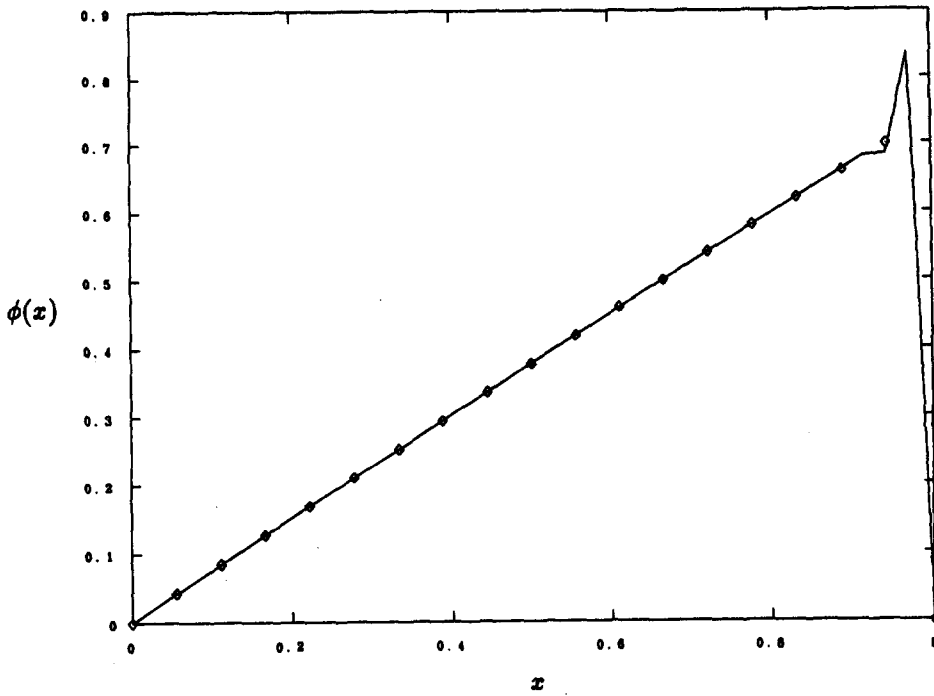


Figure 2. At $t = 1$, the accurate solution (dots) vs. LS solution of $h = 1/18$, $k = 1/18$ (line).

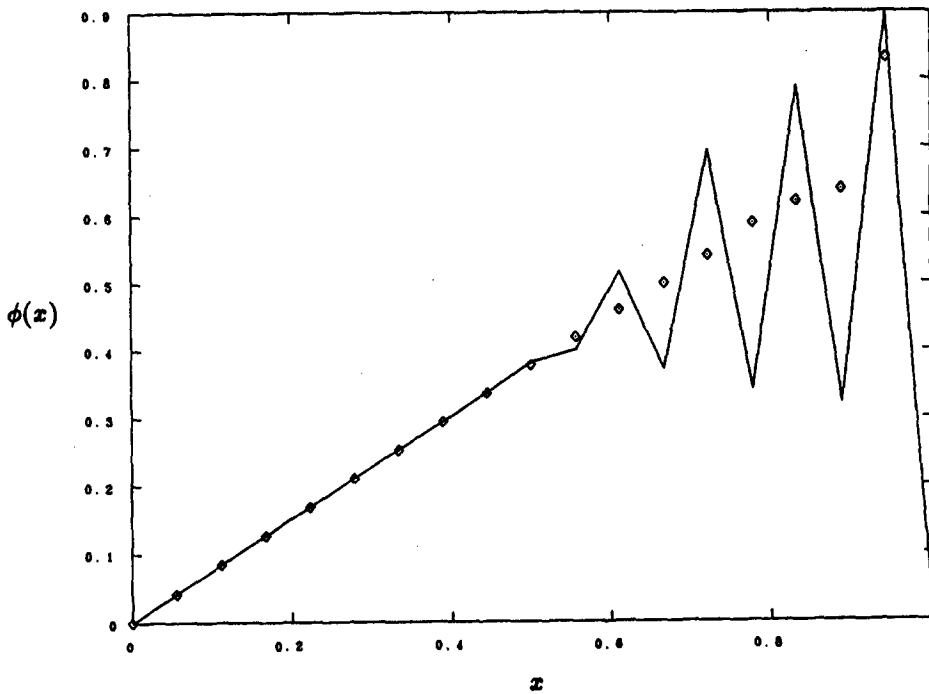


Figure 3. At $t = 1$, with $h = 1/18$, LS solution of $k = 1/18$ (dots) vs. $k = 0.01$ (line).

PROOF. Since (4.10) is valid for $(\eta, e) \in S_1 \times V_0$, it follows that

$$\|e - \eta\beta\| \leq \left(\frac{2\varepsilon}{k}\right)^{1/2} \|(\eta, e)\|_b \leq \left(\frac{2\varepsilon\delta}{C_2}\right)^{1/2} \|(\eta, e)\|_b.$$

Therefore, using triangle inequality and (4.8), we have

$$\|e\| \leq \|e - \eta\beta\| + |\beta| \|\eta\| \leq \left\{ \left(\frac{2\varepsilon\delta}{C_2}\right)^{1/2} + |\beta| \right\} \|(\eta, e)\|_b.$$

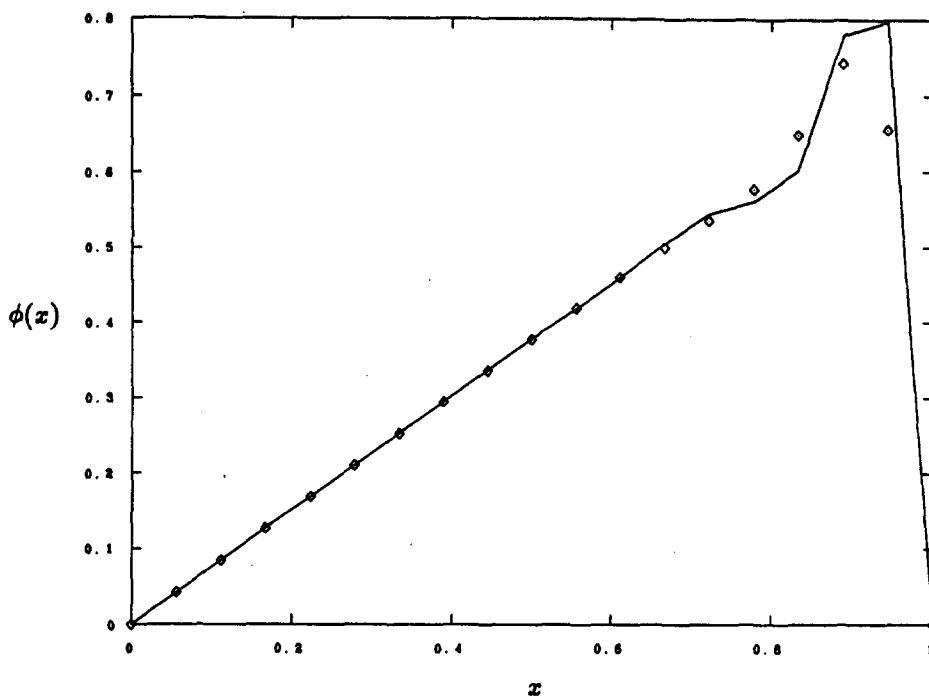


Figure 4. At $t = 1$, with $h = 1/18$, WDL solution of $k = 1/18$ (dots) vs. $k = 0.001$ (line).

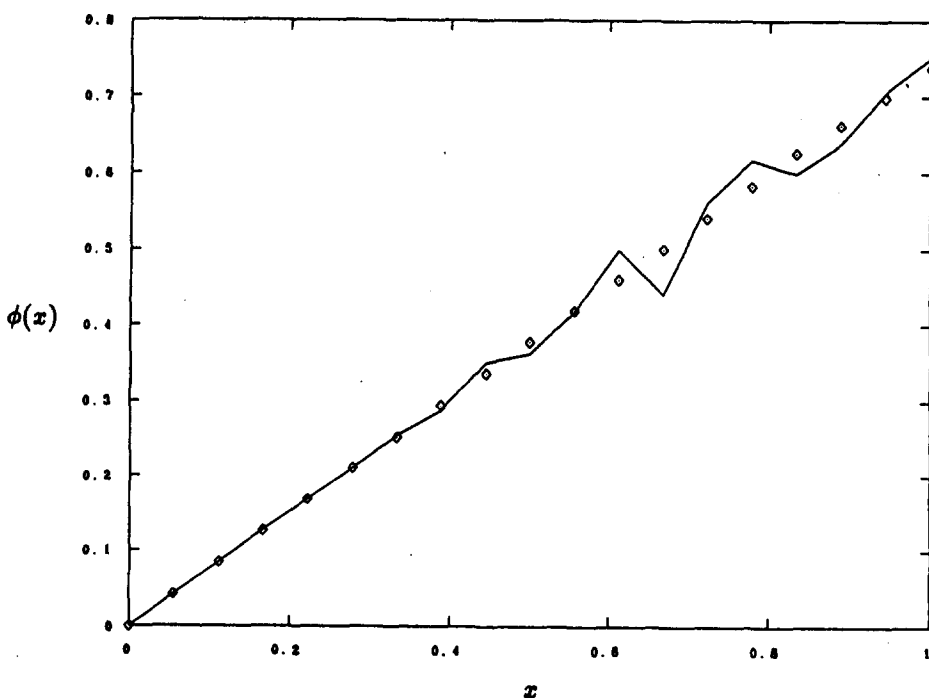


Figure 5. At $t = 1$, with $h = 1/18$, LS (5.2) solution of $k = 1/18$ (dots) vs. $k = 0.001$ (line).

Note that the proofs of (WDL) are much simpler than those of (LS). This is because that we use the condition $k = O(h)$ in the proofs all the time. This assumption is natural since we are approximating the semidiscrete problem (1.4) and the solution of (1.4) is only $O(k^2)$ as indicated in (1.5). Note that from the proofs, these analysis can also be carried out for the purely implicit problem (1.3).

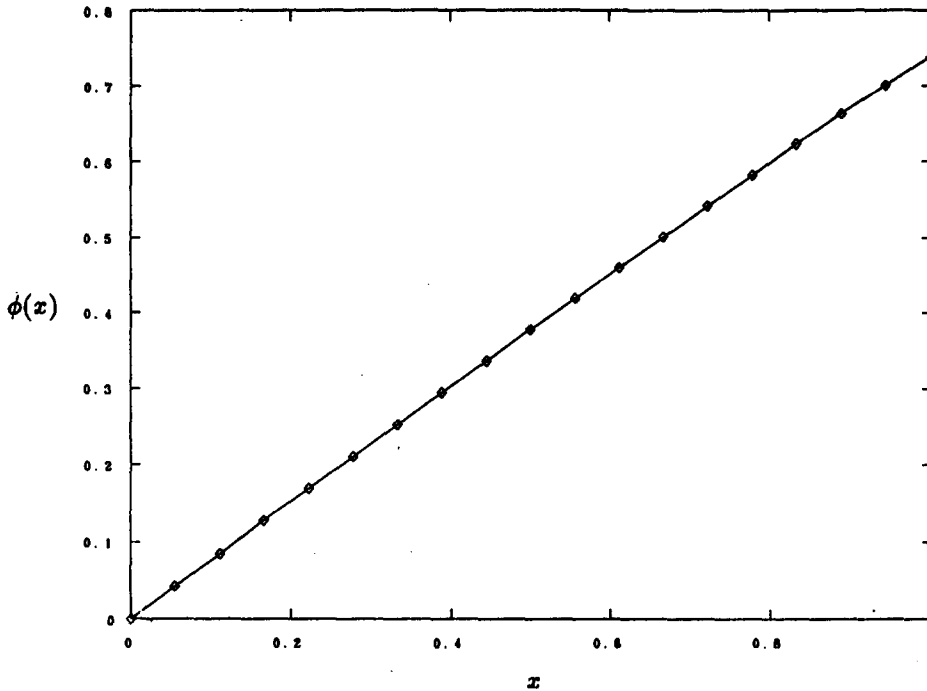


Figure 6. At $t = 1$, with $h = 1/18$, **WDLs** (5.3) solution of $k = 1/18$ (dots) vs. $k = 0.001$ (line).

The least squares (**LS**) method and weighted least squares (**WDLs**) method differ in many aspects. First, the approximation in $|||(\eta, e)|||_b$ is optimal while in $|||(\eta, e)|||$ is only suboptimal. Moreover, when the (**WDLs**) method is used, using (4.10) and argument similar to that of Theorem 4.3, the approximations in

$$\|e - \eta\beta + \varepsilon \nabla \eta\|, \quad \left| n + \frac{1}{2} k \operatorname{div} e \right|$$

are both optimal. Note that even (4.8) depends on ε , it is better than that of the (**LS**) method. One crucial difference is that (4.9) shows that the error in $|\operatorname{div} e|$ is optimal and independent of ε . In the case that ε is close to zero, this estimate is essential since it gives the approximation for $\eta\beta$. From the proofs, these are obtained without using the solvability of the dual problem (3.7), the extra regularity (3.8), or the Grid Decomposition Property (2.10). Therefore the (**WDLs**) method is far superior than the (**LS**) method and should be used in the approximation, especially when $\varepsilon < h$. In fact, in the case of $\varepsilon < h$, we can choose $\delta = \varepsilon$ and then the estimate (4.9) and (4.11) will be independent of ε .

5. NUMERICAL RESULTS

In this section we report the results of computations which illustrate the (**LS**) and (**WDLs**) methods. Although the analysis in Section 3 and 4 deal only with linear equations, the computations will be performed on the solution of a nonlinear equation. In the following, the numerical experiments deal with the Burgers' equation with $\varepsilon = 0.0001$:

$$\begin{aligned} \phi_t + \phi\phi_x - \varepsilon\phi_{xx} &= 0, & \text{in } \Omega \times I &= (0, 1) \times (0, \infty), \\ \phi(x, 0) &= \sin(\pi x), & 0 < x < 1, \\ \phi(0, t) = \phi(1, t) &= 0, & t \in I. \end{aligned} \tag{5.1}$$

Note that there is no analytical solution for this problem. Therefore it is not possible to estimate the L_2 errors occurred in the approximations. However, we will compare the results with the accurate solutions obtained by Christie and Mitchell [1].

In all the computations, linear elements are used for both S_h and V_h and the solution at $t = 1$ is calculated. Based on the (LS) and (WDLs) with $h = 1/18$ and $k = 1/18$, the numerical results obtained at $t = 1.0$ are shown in Table 1. The quoted accurate solution was computed by a Petrov-Galerkin method with fully upwinded cubic functions and a very small value of h [1]. Observe that, except the point next to the boundary, there is a great match between the result of (LS) and the accurate solution. Indeed, this result can be improved by reducing h and k simultaneously. Figure 2 shows the comparison between the accurate solution and the (LS) solution for $h = 1/36$ and $k = 1/36$. We further investigate the convergence of the (LS) method by fixing $h = 1/18$ and letting k go to zero. Computations were performed for $k = 0.01$ and $k = 0.001$. The comparative results of $k = 1/18$ and $k = 0.01$ are presented in Figure 3. Observe that large oscillation occurs when $k = 0.01$. In fact, as k continues to decrease with h fixed, the oscillation becomes more severe as it is also observed in the case when $k = 0.001$. This is similar to those obtained from the Galerkin method. As for the (WDLs) formulation, the results in Table 1 is inferior to that of the (LS). However, when h is fixed and k decreases, the oscillation is only restricted to few points next to the right hand boundary $x = 1$. This is demonstrated in Figure 4 where results of $k = 1/18$ and $k = 0.001$, with $h = 1/18$, are compared and the results are almost identical. These show that, in the case when ε is small, the method (WDLs) is stable while (LS) is not stable.

Table 1. At $t = 1$, accurate solution vs. LS and WDLs solutions obtained using $h = 1/18$, and $k = 1/18$.

x	Accurate solution	LS solution	WDLs solution
.00000E+00	.00000E+00	.00000E+00	.00000E+00
.55556E-01	.42200E-01	.42125E-01	.42128E-01
.11111E+00	.84300E-01	.84235E-01	.84241E-01
.16667E+00	.12630E+00	.12631E+00	.12632E+00
.22222E+00	.16840E+00	.16835E+00	.16836E+00
.27778E+00	.21030E+00	.21032E+00	.21033E+00
.33333E+00	.25220E+00	.25220E+00	.25222E+00
.38889E+00	.29390E+00	.29398E+00	.29398E+00
.44444E+00	.33550E+00	.33566E+00	.33563E+00
.50000E+00	.37690E+00	.37728E+00	.37728E+00
.55556E+00	.41820E+00	.41879E+00	.41905E+00
.61111E+00	.45920E+00	.45960E+00	.46051E+00
.66667E+00	.50000E+00	.49925E+00	.49984E+00
.72222E+00	.54040E+00	.54086E+00	.53587E+00
.77778E+00	.58060E+00	.58836E+00	.57730E+00
.83333E+00	.62030E+00	.62087E+00	.64938E+00
.88889E+00	.65960E+00	.63898E+00	.74272E+00
.94444E+00	.69830E+00	.83053E+00	.65594E+00
.10000E+01	.00000E+00	.00000E+00	.00000E+00

In the above computations, both (LS) and (WDLs) give bad approximations to the point next to $x = 1$. This is due to the boundary layer effect. In fact, when ε is zero, the problem (5.1) is reduced to a hyperbolic problem and the right hand boundary (the outflow boundary) condition is no longer valid. However, when ε is small, the boundary condition at $x = 1$ is still necessary to ensure the uniqueness of the solution of (5.1). This leads us to consider the following least

Table 2. At $t = 1$, with $h = 1/18$, LS (5.2) solution vs. WDLS (5.3) solutions obtained using various k .

x	LS(5.2) $k = 1/18$	WDLS(5.3) $k = 1/18$	WDLS(5.3) $k = 1/9$	WDLS(5.3) $k = 1/36$
.00000E+00	.00000E+00	.00000E+00	.00000E+00	.00000E+00
.55556E-01	.42125E-01	.42128E-01	.41775E-01	.42212E-01
.11111E+00	.84235E-01	.84241E-01	.83543E-01	.84406E-01
.16667E+00	.12631E+00	.12632E+00	.12529E+00	.12657E+00
.22222E+00	.16835E+00	.16836E+00	.16702E+00	.16868E+00
.27778E+00	.21032E+00	.21033E+00	.20871E+00	.21072E+00
.33333E+00	.25220E+00	.25222E+00	.25035E+00	.25267E+00
.38889E+00	.29399E+00	.29402E+00	.29195E+00	.29451E+00
.44444E+00	.33566E+00	.33569E+00	.33347E+00	.33622E+00
.50000E+00	.37724E+00	.37722E+00	.37493E+00	.37776E+00
.55556E+00	.41871E+00	.41860E+00	.41630E+00	.41909E+00
.61111E+00	.45980E+00	.45983E+00	.45756E+00	.46024E+00
.66667E+00	.50008E+00	.50090E+00	.49869E+00	.50125E+00
.72222E+00	.54027E+00	.54175E+00	.53963E+00	.54214E+00
.77778E+00	.58287E+00	.58222E+00	.58033E+00	.58280E+00
.83333E+00	.62620E+00	.62210E+00	.62073E+00	.62285E+00
.88889E+00	.66206E+00	.66125E+00	.66090E+00	.66189E+00
.94444E+00	.69852E+00	.69993E+00	.70109E+00	.70002E+00
.10000E+01	.73722E+00	.73880E+00	.74167E+00	.73819E+00

squares formulations. To be precise, we let the boundary $\Gamma = \Gamma_- \cup \Gamma_+$, where Γ_- is the inflow boundary and Γ_+ is the outflow boundary. Therefore, instead of minimizing (2.6) or (4.1) over the spaces S_1 and V_0 , we minimize

$$\|u - \phi\beta + \varepsilon \nabla \phi\|^2 + \left| \phi + \frac{1}{2} k \operatorname{div} u - g \right|^2 + h \int_{\Gamma_+} |\phi|^2, \quad (5.2)$$

or

$$\frac{kh}{2\varepsilon\delta} \|u - \phi\beta + \varepsilon \nabla \phi\|^2 + \left| \phi + \frac{1}{2} \operatorname{div} u - g \right|^2 + h \int_{\Gamma_+} |\phi|^2, \quad (5.3)$$

over $(\phi, u) \in S_1^- \times V_0$, where

$$S_1^- = \{\psi \mid \psi \in H^1(\Omega), \psi = 0 \text{ on } \Gamma_-\}.$$

Note that (5.2) and (5.3) are the (LS) and (WDLS) with weakly imposed boundary condition at $x = 1$, respectively. Using (5.2) and (5.3), with $h = 1/18$ and $k = 1/18$, the results are presented in Table 2. Observe that excellent agreement of the results with that of the accurate solution except at $x = 1$. The inaccuracy at $x = 1$ is not serious since the boundary condition at $x = 1$ is known. Table 2 also includes the results based on (5.3) with $h = 1/18$, $k = 1/9$ and $h = 1/18$, $k = 1/36$. These illustrate that the method based on (5.3) is stable independent of the ratio between h and k . Even though the result based on (5.2) seems to give great result in Table 2, it will again produce oscillatory solutions when h is fixed and k decreases to zero. The

Table 3. At $t = 1$, **WDLS(5.3)** solutions using $h = k = 1/36$ and $h = k = 1/72$.

x	WDLS(5.3)	
	$h = k = 1/36$	$h = k = 1/72$
.00000E+00	.00000E+00	.00000E+00
.55556E-01	.42133E-01	.42135E-01
.11111E+00	.84248E+01	.84252E-01
.16667E+00	.12633E+00	.12633E+00
.22222E+00	.16836E+00	.16836E+00
.27778E+00	.21031E+00	.21031E+00
.33333E+00	.25217E+00	.25217E+00
.38889E+00	.29393E+00	.29391E+00
.44444E+00	.33554E+00	.33552E+00
.50000E+00	.37700E+00	.37696E+00
.55556E+00	.41828E+00	.41821E+00
.61111E+00	.45934E+00	.45923E+00
.66667E+00	.50015E+00	.50001E+00
.72222E+00	.54068E+00	.54049E+00
.77778E+00	.58088E+00	.58063E+00
.83333E+00	.62071E+00	.62038E+00
.88889E+00	.66010E+00	.65968E+00
.94444E+00	.69894E+00	.69846E+00
.10000E+01	.73744E+00	.73677E+00

results based on (5.2) are reported in Figure 5. Note that the oscillation is not as severe as that of Figure 3. In Figure 6, these results are obtained based on (5.3) with $h = 1/18$, $k = 0.01$ and $k = 0.001$. No oscillation is observed in this case. We remark here that the formulation (5.3) based on the purely implicit scheme (1.3) also gives convergence results. Finally we include Table 3 to illustrate the fast convergence of the method based on (5.3).

Concluding from the above, in the case when ε is small, the method based on (5.3) is far superior than those obtained from the (LS), (WDLS) and (5.2). Moreover, both the methods (WDLS) and (5.3) give convergence results independent of the ratio of h and k . Note that from the theory, the condition that $k = O(h)$ is necessary to ensure the optimal convergence. However, the numerical results based on (5.3) and (WDLS) indicate that even in the case $h = 1/18$, $k = 0.001$, here k is worse than $O(h^2)$, we still have convergence as shown in Figure 4 and 6. Thus, the condition $k = O(h)$ is not a limitation. Moreover, as illustrated in the above, the results obtained using these least squares formulations show great efficiency as opposed to the Petrov-Galerkin method used in [1]. Above all, the upwinding is build in all these formulations in a natural way and requires no "free" parameter as often needed in existing upwinding methods [1-5].

REFERENCES

1. I. Christie and A.R. Mitchell, Upwinding of high order Galerkin Methods in conduction-convection problems, *Int. J. Num. Meth. Eng.* 12, 1764-1771 (1978).
2. D.F. Griffiths and A.R. Mitchell, On generating upwind finite element methods, In *Finite Element Methods for Convection-Dominated Flows*, (Edited by T.J.R. Hughes), pp. 91-104, ASME Monograph AMD-34, (1979).
3. C. Johnson, Streamline diffusion methods for problems in fluid mechanics, In *Finite Elements in Fluids*, (Edited by R.H. Gallagher *et al.*), Vol. 6, Wiley, New York, (1985).

4. C. Johnson and U. Nävert, An analysis of some finite element methods for convection-diffusion problems, In *Analytical and Numerical Approaches to Asymptotic Problems*, (Edited by O. Axelsson et al.), pp. 99–116, North-Holland, Amsterdam, (1981).
5. C. Johnson, U. Nävert and J. Pitkaranta, Finite element methods for linear hyperbolic problems, *Comp. Methods Appl. Mech. Eng.* **45**, 285–312 (1984).
6. J.H. Bramble and V. Thomée, Semidiscrete-least squares methods for a parabolic boundary value problem, *Math. Comp.* **26**, 633–648 (1972).
7. T.F. Chen, On least squares approximations to compressible flow problems, *Numerical Methods for Partial Differential Equations* **2**, 207–228 (1986).
8. T.F. Chen and G.J. Fix, Least squares finite element simulation of transonic flows, *Applied Numerical Mathematics* **2**, 399–408 (1986).
9. G.J. Fix, M.D. Gunzburger and R.A. Nicolaides, On mixed finite element methods for first order elliptic systems, *Numer. Math.* **37**, 29–48 (1981).
10. G. Strang and G.J. Fix, *An Analysis of the Finite Element Method*, Prentice-Hall, Inc., Englewood Cliffs, NJ, (1973).
11. J.L. Lions and E. Magenes, *Nonhomogeneous Boundary Value Problems*, Springer-Verlag, New York, (1973).