

First-Order System Least Squares for Second-Order Partial Differential Equations: Part I



Z. Cai; R. Lazarov; T. A. Manteuffel; S. F. McCormick

SIAM Journal on Numerical Analysis, Vol. 31, No. 6 (Dec., 1994), 1785-1799.

Stable URL:

<http://links.jstor.org/sici?sici=0036-1429%28199412%2931%3A6%3C1785%3AFSLSFS%3E2.0.CO%3B2-0>

SIAM Journal on Numerical Analysis is currently published by Society for Industrial and Applied Mathematics.

Your use of the JSTOR archive indicates your acceptance of JSTOR's Terms and Conditions of Use, available at <http://www.jstor.org/about/terms.html>. JSTOR's Terms and Conditions of Use provides, in part, that unless you have obtained prior permission, you may not download an entire issue of a journal or multiple copies of articles, and you may use content in the JSTOR archive only for your personal, non-commercial use.

Please contact the publisher regarding any further use of this work. Publisher contact information may be obtained at <http://www.jstor.org/journals/siam.html>.

Each copy of any part of a JSTOR transmission must contain the same copyright notice that appears on the screen or printed page of such transmission.

For more information on JSTOR contact jstor-info@umich.edu.

©2003 JSTOR

FIRST-ORDER SYSTEM LEAST SQUARES FOR SECOND-ORDER PARTIAL DIFFERENTIAL EQUATIONS: PART I*

Z. CAI[†], R. LAZAROV[‡], T. A. MANTEUFFEL[§], AND S. F. MCCORMICK[§]

Dedicated to Seymour Parter on the occasion of his 65th birthday.

Abstract. This paper develops ellipticity estimates and discretization error bounds for elliptic equations (with lower-order terms) that are reformulated as a least-squares problem for an equivalent first-order system. The main result is the proof of ellipticity, which is used in a companion paper to establish optimal convergence of multiplicative and additive solvers of the discrete systems.

Key words. least-squares discretization, second-order elliptic problems, Rayleigh–Ritz, finite elements

AMS subject classifications. 65N30

1. Introduction. The purpose of this paper is to analyze the least-squares finite element method for second-order convection-diffusion equations written as a first-order system. In general, the standard Galerkin finite element methods applied to non-self-adjoint elliptic equations with significant convection terms exhibit a variety of deficiencies, including oscillations or nonmonotonicity of the solution and poor approximation of its derivatives. A variety of stabilization techniques, such as up-winding, Petrov–Galerkin, and stream-line diffusion approximations, have been introduced to eliminate these and other drawbacks of standard Galerkin methods. Yet, although significant progress has been made, convection-diffusion problems remain among the more difficult problems to solve numerically.

Substantial progress in overcoming the deficiencies of standard Galerkin methods has been made by employing least-squares principles. From a formal point of view, the least-squares method for the linear operator equation $Lp = f$ leads to an approximation of the normal equation $L^*Lp = L^*f$, where L^* is the adjoint of L in the inner product generated by the least-squares norm. Obviously, L^*L is self-adjoint and nonnegative. However, this approach loses its advantages when directly applied to approximating second-order elliptic problems, because of the resulting fourth-order system that requires much more smoothness of the solution, thereby preventing direct use of standard finite element spaces and effectively squaring the condition number of the discrete operator.

In the last thirty years, considerable attention has been directed towards developing alternative methods by introducing physically meaningful new dependent variables (fluxes, velocity, strains and stresses, etc.) that transform the corresponding second-order elliptic problem into a system of equations of first order. The resulting system can then be posed in a weak sense and approximated by finite element methods. In most cases, this procedure leads to a saddle point problem. Due largely to Babuška [2]

* Received by the editors December 8, 1993; accepted for publication (in revised form) March 29, 1994. This work was sponsored by Air Force Office of Scientific Research grant AFOSR-86-0126 and by National Science Foundation grant DMS-8704169.

[†] Center for Applied Mathematical Sciences, Department of Mathematics, University of Southern California, 1042 W. 36th Place, DRB 155, Los Angeles, California 90089-1113 (zcaiw@uhsd.usc.edu).

[‡] Department of Mathematics, Texas A&M University, College Station, Texas 77843-3368, and Institute of Mathematics, Bulgarian Academy of Sciences, Sofia, Bulgaria (lazarov@math.tamu.edu).

[§] Program in Applied Mathematics, Campus Box 526, University of Colorado at Boulder, Boulder, Colorado 80309-0526 (tmanteuf@boulder.colorado.edu and stevem@boulder.colorado.edu).

and Brezzi [5], it is now well understood that the finite element spaces approximating different physical quantities (pressure and velocity, or temperature and flux, or displacement and stresses, etc.) cannot be chosen independently if one wants to have an unconditionally stable scheme of optimal approximation order. In particular, these spaces are usually chosen to satisfy the so-called inf-sup condition of Ladyzhenskaya, Babuška, and Brezzi [23], [2], [5].

A general theory of the least-squares method for approximating elliptic boundary value problems of Agmon–Douglas–Nirenberg (ADN) type has been developed by Aziz, Kellogg, and Stephens in [1]. The method involves the minimization of a least-squares functional that consists of a weighted sum of the residuals occurring in the equations and the boundary conditions. The weights occurring in the least-squares functional are determined by the indices that enter into the definition of the ADN boundary value problem. This approach generalizes both the least-squares method of Jespersen [20], which is for systems arising from reformulating the Poisson equation by introducing the gradient of the solution as a new unknown, and the method of Wendland [28], which is for elliptic systems of Cauchy–Riemann type. However, this approach leads to an algebraic problem with condition number $O(h^{-4})$, where h is the grid size.

In a recent work by Bochev and Gunzburger [4], the ADN approach was extended to a velocity–vorticity–pressure formulation of Stokes flows with rigorous error analysis and optimal convergence rates. Further extension of these ideas to Navier–Stokes equations with interesting numerical experiments for two-dimensional (2D) problems appears in [3].

The most important applications of the least-squares method are connected with continuum mechanics (cf. [14], [21], [3], and [24]). A common feature of all approaches is the introduction of new dependent variables: pressure in elastic problems (see [14]) or vorticity in flow problems (see [3]), for example. This has the effect of transforming the second-order elliptic problem into a system of first order. In most cases, the resulting system will not be of Petrovsky type (cf. [28]). On the other hand, one can add to this system a variety of compatibility conditions to overcome this limitation. For example, the system $-\operatorname{div} \mathbf{u} = f$, $\mathbf{u} = \mathbf{grad} p$ that decomposes the Poisson equation $-\Delta p = f$ can be effectively augmented by the compatibility equation $\operatorname{curl} \mathbf{u} = \mathbf{0}$ (see [22] and [24]). Another possibility, which was explored by Chen and Fix in [10] and [11] for fluid flow, is to eliminate the variable p altogether in favor of a least-squares functional based on $\operatorname{curl} \mathbf{u} = \mathbf{0}$ and $\operatorname{div} \mathbf{u} + f = 0$.

The least-squares approach represents a general methodology that can produce a variety of algorithms, depending on such choices as the least-squares norm, the least-squares system, and the boundary treatment, and that can lead to formulations that have substantially different properties. For example, some formulations lead to self-adjoint problems, others yield optimal errors, others have unconditional stability, and others lead to optimal conditioning of the resulting discrete system. Among the deficiencies of some of the incarnations of this approach are non-self-adjointness, conditional stability, $O(h^{-4})$ condition number of the resulting discrete problem, and special essential boundary conditions. The purpose here is to develop a least-squares approach that does not exhibit these limitations. Another deficiency of conventional least-square methodology is the general lack of an efficient solver for the resulting discrete system. This is the subject of a companion paper [8].

Our attention here to least squares was motivated by the numerical experiments of Chen and Fix in [11] and of Carey and Shen in [9], where a second-order equation

is solved using least-squares finite elements. In particular, the numerical experiments from [9] show that the finite element approximations of the unknown function p and its gradient \mathbf{u} converge at rates depending only on the approximation properties of the finite element spaces of piecewise polynomials of degree k and r for p and \mathbf{u} , respectively. This is an interesting phenomenon since least squares leads to a coupled system for p and \mathbf{u} . In [26], a theory that explains some of these computational results was developed. These ideas are extended in [25] to multidimensional self-adjoint equations of second-order split into a system of first order. The main point in [25] is that the least-squares functional generates a bilinear form that is continuous and elliptic in a properly defined subspace of $H(\text{div}) \times H^1$ and, therefore, standard finite element theory can be applied. In particular, one can use Raviart–Thomas mixed finite element spaces from [27] in order to approximate $H(\text{div})$. However, any other finite-dimensional subspace of $H(\text{div})$ can be used, since the approximating space need not satisfy the inf-sup condition.

In this paper, our goal is to extend the theory in [25] to systems arising from splitting convection-diffusion and reaction-diffusion equations into a system of equations of first order. The main result here is in §3, where ellipticity with respect to the $H(\text{div}) \times H^1$ norm of the bilinear form corresponding to the least-squares functional is established under very general assumptions on the differential operator. Other possible forms of the least-squares functional are suggested in §4. Approximation by the finite element method and its error analysis is then carried out in §5 in a straightforward manner, showing that the error in the $H(\text{div}) \times H^1$ norm is optimal. Moreover, we show in §6 that the condition number for the resulting linear system is at most $O(h^{-2})$.

In the engineering literature, there is an impressive number of papers on least squares for a variety of applied problems. For a review of these applications up to the mid-1970's, see [13].

2. Problem formulation. Let Ω be a bounded domain in \mathbf{R}^n , $n = 2$ or 3 , with Lipschitz boundary $\partial\Omega$. We consider the boundary value problem

$$(2.1) \quad -\text{div } A \mathbf{grad} p + Xp = f \quad \text{in } \Omega,$$

$$(2.2) \quad p = 0 \quad \text{on } \Gamma_D,$$

$$(2.3) \quad \mathbf{n} \cdot A \mathbf{grad} p = 0 \quad \text{on } \Gamma_N,$$

for $f \in L^2(\Omega)$, where $A(x)$ is an $n \times n$ symmetric matrix of functions in $L^2(\Omega)$, X is a linear differential operator of order at most 1, $\partial\Omega = \Gamma_D \cup \Gamma_N$ is the boundary of Ω , and \mathbf{n} is the outward unit vector normal to the boundary. We assume that A is uniformly positive definite and scaled appropriately; that is, there exist positive constants

$$(2.4) \quad 0 < \lambda \leq 1 \leq \Lambda$$

such that

$$(2.5) \quad \lambda \boldsymbol{\xi}^T \boldsymbol{\xi} \leq \boldsymbol{\xi}^T A \boldsymbol{\xi} \leq \Lambda \boldsymbol{\xi}^T \boldsymbol{\xi}$$

for all $\boldsymbol{\xi} \in \mathbf{R}^n$ and almost all $x \in \overline{\Omega}$.

Possible choices for X include

$$(2.6) \quad \begin{aligned} (a) \quad & Xp = \operatorname{div}(\mathbf{b}p), \quad \mathbf{b} = (b_1(x), \dots, b_n(x)) \in (L^2(\Omega))^n; \\ (b) \quad & Xp = \mathbf{a} \cdot \mathbf{grad} p + cp, \quad \mathbf{a} \in (L^2(\Omega))^n, \quad c(x) \in L^2(\Omega). \end{aligned}$$

In particular, a real Helmholtz problem arises with $\mathbf{a} \equiv \mathbf{0}$ and $c(x) = -k^2$ in (2.6(b)).

The classical Sobolev spaces $H^m(\Omega)$, with m th norm $\|\cdot\|_{m,\Omega}$ and seminorms $|\cdot|_{i,\Omega}$, $0 \leq i \leq m$, are employed and, as usual, $L^2(\Omega) = H^0(\Omega)$. Denote the corresponding norms on product space $(H^m(\Omega))^n$ by $\|\cdot\|_{m,\Omega,n}$ and $|\cdot|_{i,\Omega,n}$. We also need

$$H(\operatorname{div}) \equiv \{\mathbf{v} \in (L^2(\Omega))^n : \operatorname{div} \mathbf{v} \in L^2(\Omega)\}.$$

We will be interested in the spaces (see (2.2) and (2.3))

$$(2.7) \quad \begin{aligned} \mathbf{W} &\equiv \{\mathbf{v} \in H(\operatorname{div}) : \mathbf{n} \cdot \mathbf{v} = 0 \text{ on } \Gamma_N\}, \\ V &\equiv \{q \in H^1(\Omega) : q = 0 \text{ on } \Gamma_D\}, \end{aligned}$$

with respective norms

$$(2.8) \quad \begin{aligned} \|\mathbf{v}\|_{H(\operatorname{div})}^2 &= \|\mathbf{v}\|_{0,\Omega,n}^2 + \|\operatorname{div} \mathbf{v}\|_{0,\Omega}^2 \equiv \|\mathbf{v}\|_{\mathbf{W}}^2, \\ \|q\|_{1,\Omega}^2 &= \|q\|_{0,\Omega}^2 + \|\mathbf{grad} q\|_{0,\Omega,n}^2 \equiv \|q\|_V^2. \end{aligned}$$

As usual, the inner product on $L^2(\Omega)$ is denoted by (\cdot, \cdot) ; i.e., for any $p, q \in L^2(\Omega)$,

$$(2.9) \quad (p, q) = \int_{\Omega} p q \, dx.$$

Likewise, the inner product on $(L^2(\Omega))^n$ is denoted by

$$(2.10) \quad (\mathbf{u}, \mathbf{v})_n = \sum_{i=1}^n \int_{\Omega} u_i v_i \, dx \quad \forall \mathbf{u}, \mathbf{v} \in (L^2(\Omega))^n.$$

Given a linear operator $X : H^1 \rightarrow L^2$, denote by $X^* : H^1 \rightarrow L^2$ its formal L^2 adjoint defined by

$$(2.11) \quad (Xp, q) = (p, X^*q)$$

for p and $q \in C_0^\infty(\Omega)$. In other words, we make the definition without regard to boundary conditions. Likewise, given a linear operator $X : H^1 \rightarrow (L^2)^n$, denote by $X^* : (H^1)^n \rightarrow L^2$ its formal L^2 adjoint defined by

$$(2.12) \quad (Xq, \mathbf{v})_n = (q, X^*\mathbf{v})$$

for all $q \in C_0^\infty(\Omega)$ and $\mathbf{v} \in (C_0^\infty(\Omega))^n$. Define the operator $\nabla : H^1 \rightarrow (L^2)^n$ by

$$(2.13) \quad \nabla q \equiv \mathbf{grad} q \equiv \left(\frac{\partial q}{\partial x_1}, \frac{\partial q}{\partial x_2}, \dots, \frac{\partial q}{\partial x_n} \right).$$

Its formal adjoint $\nabla^* : (H^1)^n \rightarrow L^2$ is then defined as

$$(2.14) \quad \nabla^*\mathbf{v} = -\operatorname{div} \mathbf{v} = -\left(\frac{\partial v_1}{\partial x_1} + \dots + \frac{\partial v_n}{\partial x_n} \right).$$

Notice that if we restrict ∇ to V and ∇^* to \mathbf{W} , then ∇^* is, in fact, the $L^2(\Omega)$ adjoint of ∇ .

Consider the least-squares functional

$$(2.15) \quad G(\mathbf{u}, p; f) = \|A\nabla p - \mathbf{u}\|_{0,\Omega,n}^2 + \|\nabla^* \mathbf{u} + Xp - f\|_{0,\Omega}^2,$$

and the associated bilinear form

$$(2.16) \quad \mathcal{F}(\mathbf{u}, p; \mathbf{v}, q) = \left(\begin{pmatrix} -I & A\nabla \\ \nabla^* & X \end{pmatrix} \begin{pmatrix} \mathbf{u} \\ p \end{pmatrix}, \begin{pmatrix} -I & A\nabla \\ \nabla^* & X \end{pmatrix} \begin{pmatrix} \mathbf{v} \\ q \end{pmatrix} \right)_{n+1}.$$

Assuming that (2.1)–(2.3) is uniquely solvable in H^1 , it is easy to see that the (unique) minimum of $G(\mathbf{u}, p; f)$ over $(\mathbf{u}, p) \in \mathbf{W} \times V$, for fixed $f \in L^2$, corresponds to the solution of (2.1)–(2.3) in the sense that p solves these equations and $\mathbf{u} = A\nabla p$. The main result of this paper will be to show that $G(\mathbf{u}, p; f)$ is elliptic with respect to the $H(\text{div}) \times H^1$ norm, that is, there exist positive constants α and β such that

$$(2.17) \quad \alpha \left(\|\mathbf{u}\|_{H(\text{div})}^2 + \|p\|_{1,\Omega}^2 \right) \leq G(\mathbf{u}, p; 0) \leq \beta \left(\|\mathbf{u}\|_{H(\text{div})}^2 + \|p\|_{1,\Omega}^2 \right)$$

for every $\mathbf{u} \in \mathbf{W}$ and $p \in V$.

The uniform bounds on A (see (2.5)) imply that the functional in (2.15) is equivalent to the functional

$$(2.18) \quad \widehat{G}(\mathbf{u}, p; f) = \|A^{1/2}\nabla p - A^{-1/2}\mathbf{u}\|_{0,\Omega,n}^2 + \|\nabla^* \mathbf{u} + Xp - f\|_{0,\Omega}^2$$

in the sense that

$$(2.19) \quad \lambda \widehat{G}(\mathbf{u}, p; f) \leq G(\mathbf{u}, p; f) \leq \Lambda \widehat{G}(\mathbf{u}, p; f)$$

for every $\mathbf{u} \in \mathbf{W}$, $p \in V$, and $f \in L^2(\Omega)$. In our proof we will be working with $\widehat{G}(\mathbf{u}, p; f)$, which can be written as

$$\widehat{G}(\mathbf{u}, p; f) = \left(\begin{pmatrix} -A^{-1/2} & A^{1/2}\nabla \\ \nabla^* & X \end{pmatrix} \begin{pmatrix} \mathbf{u} \\ p \end{pmatrix} - \begin{pmatrix} \mathbf{0} \\ f \end{pmatrix}, \begin{pmatrix} -A^{-1/2} & A^{1/2}\nabla \\ \nabla^* & X \end{pmatrix} \begin{pmatrix} \mathbf{u} \\ p \end{pmatrix} - \begin{pmatrix} \mathbf{0} \\ f \end{pmatrix} \right)_{n+1},$$

and the associated bilinear form

$$(2.20) \quad \widehat{\mathcal{F}}(\mathbf{u}, p; \mathbf{v}, q) = \left(\begin{pmatrix} -A^{-1/2} & A^{1/2}\nabla \\ \nabla^* & X \end{pmatrix} \begin{pmatrix} \mathbf{u} \\ p \end{pmatrix}, \begin{pmatrix} -A^{-1/2} & A^{1/2}\nabla \\ \nabla^* & X \end{pmatrix} \begin{pmatrix} \mathbf{v} \\ q \end{pmatrix} \right)_{n+1}.$$

If we restrict our attention to functions with the correct properties so that we can, for the moment, ignore the issues of smoothness and boundary conditions, then this form can be written as

$$(2.21) \quad \widehat{\mathcal{F}}(\mathbf{u}, p; \mathbf{v}, q) = \left(\begin{pmatrix} A^{-1} + \nabla\nabla^* & \nabla(X - I) \\ (X^* - I)\nabla^* & \nabla^*A\nabla + X^*X \end{pmatrix} \begin{pmatrix} \mathbf{u} \\ p \end{pmatrix}, \begin{pmatrix} \mathbf{v} \\ q \end{pmatrix} \right)_{n+1}.$$

The $H(\text{div}) \times H^1$ norm is induced by the inner product

$$(2.22) \quad (\mathbf{u}, p; \mathbf{v}, q)_{H(\text{div}) \times H^1} = (\mathbf{u}, \mathbf{v})_n + (\nabla^* \mathbf{u}, \nabla^* \mathbf{v}) + (p, q) + (\nabla p, \nabla q)_n.$$

Similar in spirit to the definition of \mathcal{F} , we define an inner product

$$(2.23) \quad \mathcal{S}(\mathbf{u}, p; \mathbf{v}, q) = (A^{-1/2}\mathbf{u}, A^{-1/2}\mathbf{v})_n + (\nabla^*\mathbf{u}, \nabla^*\mathbf{v}) + (p, q) + (A^{1/2}\nabla p, A^{1/2}\nabla q)_n,$$

which satisfies equivalence bounds

$$(2.24) \quad \frac{1}{C}\mathcal{S}(\mathbf{u}, p; \mathbf{u}, p) \leq (\mathbf{u}, p; \mathbf{u}, p)_{H(\text{div}) \times H^1} \leq C\mathcal{S}(\mathbf{u}, p; \mathbf{u}, p),$$

where $C = \max\{1/\lambda, \Lambda\}$. (Recall $0 < \lambda \leq 1 \leq \Lambda$.) Again taking liberties with smoothness and boundary conditions, we can write

$$(2.25) \quad \mathcal{S}(\mathbf{u}, p; \mathbf{v}, q) = \left(\begin{pmatrix} A^{-1} + \nabla\nabla^* & 0 \\ 0 & I + \nabla^*A\nabla \end{pmatrix} \begin{pmatrix} \mathbf{u} \\ p \end{pmatrix}, \begin{pmatrix} \mathbf{v} \\ q \end{pmatrix} \right).$$

Notice the similarity between the diagonal terms of the matrices defining the bilinear forms in (2.21) and (2.25). The main result of this paper can be viewed as demonstrating the spectral equivalence of these two bilinear forms. Of course, we must pay close attention to the issues of smoothness and boundary conditions to make sure the equivalence holds on all of $\mathbf{W} \times V$.

3. Main result. The main objective of this paper is to establish equivalence of the forms \mathcal{S} in (2.23) and \mathcal{F} in (2.16) on $\mathbf{W} \times V$. If problem (2.1)–(2.3) is invertible, and $\partial\Omega$ and the coefficients of A and X are $C^{1,1}$ (first derivative is Lipschitz continuous), then the problem will be $H^2(\Omega)$ regular and the proof flows easily. However, with a little care the result can be proved under more general hypotheses for (2.1)–(2.3). In particular, we wish to include the important case in which A is discontinuous and $\partial\Omega$ is polygonal. In the next few remarks we discuss certain results that are central to the proof.

Remark 3.1. Our theory assumes that either $\Gamma_D \neq \emptyset$ or an additional constraint is imposed on V , such as $\int_{\Omega} p \, dx = 0$, so that a Poincaré–Friedrichs inequality holds; that is, there exists a constant $d > 0$ depending only on the domain Ω and the uniform bounds on A (see (2.5)) such that

$$(3.1) \quad \|p\|_{0,\Omega}^2 \leq d\|A^{1/2}\nabla p\|_{0,\Omega}^2$$

for $p \in V$.

Remark 3.2. We will also assume that (2.1)–(2.3) is invertible in $H^1(\Omega)$; that is, for any $f \in H^{-1}(\Omega)$ there is a unique $p \in V$ such that

$$(3.2) \quad (A\nabla p, \nabla v)_n + (Xp, v) = (f, v)$$

for every $v \in V$. This is true for the case $X = 0$, $X = I$, or X as given by (2.6(a) or 2.6(b)) (see [17]). If $f \in L^2(\Omega)$, then the solution p will satisfy $A\nabla p \in \mathbf{W}$ (see [15]). We will make particular use of this result with $X \equiv 0$ and $X \equiv I$. With this in mind, we make the definition

$$(3.3) \quad D = \{p \in V : A\nabla p \in \mathbf{W}\},$$

and note that invertibility of (2.1)–(2.3) means that it defines a bijection from $L^2(\Omega)$ to D .

Remark 3.3. While it is almost implicit in the assumption that (2.1)–(2.3) is invertible, we will make overt use of the bound

$$(3.4) \quad \|Xp\|_{0,\Omega} \leq \eta\|A^{1/2}\nabla p\|_{0,\Omega}$$

for some $\eta > 0$ and every $p \in D$. This bound clearly follows from (3.1) if X has the form (2.6(a)) and the coefficients are sufficiently smooth, or if X has the form (2.6(b)).

Remark 3.4. Finally, we will make use of the inequality

$$(3.5) \quad \|\nabla^* A \nabla p + Xp\|_{0,\Omega}^2 \geq \delta \|\nabla^* A \nabla p\|_{0,\Omega}^2$$

for some $\delta > 0$ and for all $p \in D$. Its validity follows from the results in [16], where it is shown that if two uniformly elliptic operators are invertible in $H^1(\Omega)$ and have the same leading part and boundary conditions, then they are $L^2(\Omega)$ -norm equivalent, even in the absence of $H^2(\Omega)$ regularity.

In addition, since $(\nabla^* A \nabla)^{-1}$ is bounded in L^2 , there is some constant K such that

$$(3.6) \quad \|p\|_{0,\Omega} \leq K \|\nabla^* A \nabla p\|_{0,\Omega}$$

for every $p \in D$. Thus, by the definition of D and integration by parts, we have

$$\|A^{1/2} \nabla p\|_{0,\Omega}^2 = (\nabla^* A \nabla p, p) \leq K (\nabla^* A \nabla p, \nabla^* A \nabla p) = K \|\nabla^* A \nabla p\|_{0,\Omega}^2$$

for all $p \in D$. Together with (3.5), this yields

$$(3.7) \quad \|\nabla^* A \nabla p + Xp\|_{0,\Omega}^2 \geq \gamma (\|\nabla^* A \nabla p\|_{0,\Omega}^2 + \|A^{1/2} \nabla p\|_{0,\Omega}^2)$$

for all $p \in D$, where $\gamma = \delta/(K+1)$.

THEOREM 3.1. *Assume that V implies the Poincaré–Friedrichs inequality (3.1), that X satisfies the bound (3.4), and that (2.1)–(2.3) is invertible in $H^1(\Omega)$. Then there exist positive constants α and β such that*

$$(3.8) \quad \mathcal{F}(\mathbf{u}, p; \mathbf{v}, q) \leq \beta \left(\|\mathbf{u}\|_{H(\text{div})}^2 + \|p\|_{1,\Omega}^2 \right)^{1/2} \left(\|\mathbf{v}\|_{H(\text{div})}^2 + \|q\|_{1,\Omega}^2 \right)^{1/2}$$

for every $\mathbf{u}, \mathbf{v} \in \mathbf{W}$ and $p, q \in V$ and

$$(3.9) \quad \mathcal{F}(\mathbf{u}, p; \mathbf{u}, p) \geq \alpha \left(\|\mathbf{u}\|_{H(\text{div})}^2 + \|p\|_{1,\Omega}^2 \right)$$

for every $\mathbf{u} \in \mathbf{W}$ and $p \in V$.

Proof. Continuity of \mathcal{F} (3.8) follows directly from assumption (2.5), definitions (2.8) and (2.16), and (3.4). The proof of the lower bound (3.9) is established in a series of steps.

Formulate an equivalent problem: By (2.19) and (2.24) it is sufficient to find a positive constant α_0 such that

$$(3.10) \quad \alpha_0 \mathcal{S}(\mathbf{u}, p; \mathbf{u}, p) \leq \widehat{\mathcal{F}}(\mathbf{u}, p; \mathbf{u}, p)$$

for all $\mathbf{u} \in \mathbf{W}$ and $p \in V$, where \mathcal{S} and $\widehat{\mathcal{F}}$ are defined in (2.23) and (2.20), respectively. Then (3.9) would follow with $\alpha = \alpha_0 \lambda \min\{\lambda, 1/\Lambda\}$.

Remove divergence-free part: For any $\mathbf{u} \in \mathbf{W}$, we may use a decomposition analogous to Weyl's theorem [15]:

$$(3.11) \quad \mathbf{u} = A \nabla \phi + \boldsymbol{\psi},$$

where $\phi \in D$ and

$$(3.12) \quad \nabla^* \boldsymbol{\psi} = 0,$$

$$(3.13) \quad \mathbf{n} \cdot \boldsymbol{\psi} = 0 \quad \text{on } \Gamma_N.$$

This is accomplished by choosing ϕ to be the weak solution of

$$(3.14) \quad \nabla^* A \nabla \phi = \nabla^* \mathbf{u},$$

$$(3.15) \quad \phi = 0 \quad \text{on } \Gamma_D,$$

$$(3.16) \quad \mathbf{n} \cdot A \nabla \phi = 0 \quad \text{on } \Gamma_N.$$

From Remark 3.2 it follows that $\phi \in D$, so $A \nabla \phi \in \mathbf{W}$. Setting $\boldsymbol{\psi} = \mathbf{u} - A \nabla \phi$, then $\boldsymbol{\psi} \in \mathbf{W}$ and (3.11)–(3.13) follows.

Equations (2.16), (2.22), and (2.23) yield

$$(3.17) \quad \mathcal{S}(\boldsymbol{\psi}, 0; \boldsymbol{\psi}, 0) = \widehat{\mathcal{F}}(\boldsymbol{\psi}, 0; \boldsymbol{\psi}, 0) = (A^{-1/2} \boldsymbol{\psi}, A^{-1/2} \boldsymbol{\psi})_n.$$

Consider the cross-product terms. Since $\boldsymbol{\psi} \in \mathbf{W}$ and $\phi, p \in V$, we can integrate by parts (see [15]) to obtain

$$(3.18) \quad \mathcal{S}(A \nabla \phi, p; \boldsymbol{\psi}, 0) = (\nabla \phi, \boldsymbol{\psi})_n = (\phi, \nabla^* \boldsymbol{\psi}) = 0,$$

$$(3.19) \quad \widehat{\mathcal{F}}(A \nabla \phi, p; \boldsymbol{\psi}, 0) = (\nabla p - \nabla \phi, \boldsymbol{\psi})_n = (p - \phi, \nabla^* \boldsymbol{\psi}) = 0.$$

This yields

$$\mathcal{S}(A \nabla \phi + \boldsymbol{\psi}, p; A \nabla \phi + \boldsymbol{\psi}, p) = \mathcal{S}(A \nabla \phi, p; A \nabla \phi, p) + \mathcal{S}(\boldsymbol{\psi}, 0; \boldsymbol{\psi}, 0),$$

$$\widehat{\mathcal{F}}(A \nabla \phi + \boldsymbol{\psi}, p; A \nabla \phi + \boldsymbol{\psi}, p) = \widehat{\mathcal{F}}(A \nabla \phi, p; A \nabla \phi, p) + \mathcal{S}(\boldsymbol{\psi}, 0; \boldsymbol{\psi}, 0).$$

Thus, it only remains to show that there exists a positive constant $\alpha_0 \leq 1$ such that

$$(3.20) \quad \alpha_0 \mathcal{S}(A \nabla \phi, p; A \nabla \phi, p) \leq \widehat{\mathcal{F}}(A \nabla \phi, p; A \nabla \phi, p)$$

for every $\phi \in D$, $p \in V$.

Define \mathcal{S}_0 and \mathcal{F}_0 and show that $c \mathcal{S}_0 \leq \mathcal{F}_0$. Let

$$\begin{aligned} \mathcal{F}_0(\phi, p; \varphi, q) &= \widehat{\mathcal{F}}(A \nabla \phi, p; A \nabla \varphi, q) \\ &= \left(\begin{pmatrix} -A^{1/2} \nabla & A^{1/2} \nabla \\ \nabla^* A \nabla & X \end{pmatrix} \begin{pmatrix} \phi \\ p \end{pmatrix}, \begin{pmatrix} -A^{1/2} \nabla & A^{1/2} \nabla \\ \nabla^* A \nabla & X \end{pmatrix} \begin{pmatrix} \varphi \\ q \end{pmatrix} \right)_{n+1} \end{aligned}$$

and

$$\mathcal{S}_0(\phi, p; \varphi, q) = \left(\begin{pmatrix} 0 & A^{1/2} \nabla \\ \nabla^* A \nabla & 0 \end{pmatrix} \begin{pmatrix} \phi \\ p \end{pmatrix}, \begin{pmatrix} 0 & A^{1/2} \nabla \\ \nabla^* A \nabla & 0 \end{pmatrix} \begin{pmatrix} \varphi \\ q \end{pmatrix} \right)_{n+1}.$$

We now show that there exists a positive constant c such that

$$(3.21) \quad c \mathcal{S}_0(\phi, p; \phi, p) \leq \mathcal{F}_0(\phi, p; \phi, p)$$

for every $\phi \in D$ and $p \in V$. We start by writing

$$(3.22) \quad \begin{pmatrix} \phi \\ p \end{pmatrix} = \begin{pmatrix} z \\ z \end{pmatrix} + \begin{pmatrix} w_1 \\ w_2 \end{pmatrix},$$

where $z, w_1 \in D$ and $w_2 \in V$ such that

$$(3.23) \quad \mathcal{S}_0(z, z; w_1, w_2) = 0,$$

and

$$(3.24) \quad \nabla^* A \nabla w_1 = -w_2.$$

This is accomplished by choosing $z, w_1 \in D$ to satisfy

$$(3.25) \quad \nabla^* A \nabla z + z = \nabla^* A \nabla \phi + p,$$

$$(3.26) \quad \nabla^* A \nabla w_1 + w_1 = \phi - p$$

(see Remark 3.2). To see this, note that adding (3.25) and (3.26) yields

$$(3.27) \quad \phi = z + w_1.$$

Since $D \subset V$, setting

$$(3.28) \quad w_2 = p - z$$

yields $w_2 \in V$ and (3.22). Substituting (3.27) into (3.26) and using (3.28) proves (3.24), which in turn yields (3.23).

Since $z \in D$, the bound (3.7) yields

$$(3.29) \quad \begin{aligned} \mathcal{F}_0(z, z; z, z) &= (\nabla^* A \nabla z + Xz, \nabla^* A \nabla z + Xz) = \|\nabla^* A \nabla z + Xz\|_{0, \Omega}^2 \\ &\geq \gamma (\|\nabla^* A \nabla z\|_{0, \Omega}^2 + \|A^{1/2} \nabla z\|_{0, \Omega}^2) = \gamma \mathcal{S}_0(z, z; z, z). \end{aligned}$$

Note that

$$(3.30) \quad \mathcal{F}_0(w_1, w_2; w_1, w_2) = \|A^{1/2} \nabla w_2 - A^{1/2} \nabla w_1\|_{0, \Omega, n}^2 + \|\nabla^* A \nabla w_1 + Xw_2\|_{0, \Omega}^2.$$

Using (3.24), the first term in (3.30) satisfies

$$(3.31) \quad \begin{aligned} &\|A^{1/2} \nabla w_2 - A^{1/2} \nabla w_1\|_{0, \Omega, n}^2 \\ &= (A \nabla w_1, \nabla w_1)_n - 2(A \nabla w_1, \nabla w_2)_n + (A \nabla w_2, \nabla w_2)_n \\ &= (A \nabla w_1, \nabla w_1)_n - 2(\nabla^* A \nabla w_1, w_2) + (A \nabla w_2, \nabla w_2)_n \\ &= (A \nabla w_1, \nabla w_1)_n + 2(\nabla^* A \nabla w_1, \nabla^* A \nabla w_1) + (A \nabla w_2, \nabla w_2)_n \\ &\geq (\nabla^* A \nabla w_1, \nabla^* A \nabla w_1) + (A \nabla w_2, \nabla w_2)_n \\ &= \mathcal{S}_0(w_1, w_2; w_1, w_2). \end{aligned}$$

Note that the integration by parts is justified because $w_1 \in D$ and $w_2 \in V$. Using the bound (3.4) and assuming, as we may, that $\eta \geq 1$, the second term in (3.30) satisfies

$$(3.32) \quad \begin{aligned} \|\nabla^* A \nabla w_1 + Xw_2\|_{0, \Omega}^2 &\leq 2(\|\nabla^* A \nabla w_1\|_{0, \Omega}^2 + \|Xw_2\|_{0, \Omega}^2) \\ &\leq 2(\|\nabla^* A \nabla w_1\|_{0, \Omega}^2 + \eta^2 \|A^{1/2} \nabla w_2\|_{0, \Omega, n}^2) \\ &\leq 2\eta^2 \mathcal{S}_0(w_1, w_2; w_1, w_2). \end{aligned}$$

Using (3.22) and the Cauchy–Schwarz inequality we have

$$\begin{aligned} \mathcal{F}_0(\phi, p; \phi, p) &= \|A^{1/2}\nabla w_2 - A^{1/2}\nabla w_1\|_{0,\Omega,n}^2 + \|\nabla^*A\nabla w_1 + Xw_2\|_{0,\Omega}^2 \\ &\quad + 2(\nabla^*A\nabla w_1 + Xw_2, \nabla^*A\nabla z + Xz) + \|\nabla^*A\nabla z + Xz\|_{0,\Omega}^2 \\ &\geq \|A^{1/2}\nabla w_2 - A^{1/2}\nabla w_1\|_{0,\Omega,n}^2 + \|\nabla^*A\nabla w_1 + Xw_2\|_{0,\Omega}^2 \\ &\quad - 2\|\nabla^*A\nabla w_1 + Xw_2\|_{0,\Omega}\|\nabla^*A\nabla z + Xz\|_{0,\Omega} + \|\nabla^*A\nabla z + Xz\|_{0,\Omega}^2 \\ &\geq \|A^{1/2}\nabla w_2 - A^{1/2}\nabla w_1\|_{0,\Omega,n}^2 + \left(1 - \frac{1}{\epsilon}\right)\|\nabla^*A\nabla w_1 + Xw_2\|_{0,\Omega}^2 \\ &\quad + (1 - \epsilon)\|\nabla^*A\nabla z + Xz\|_{0,\Omega}^2 \end{aligned}$$

for any $\epsilon > 0$. Considering only the case $\epsilon < 1$ so that $1 - 1/\epsilon < 0$, and using (3.29), (3.31), and (3.32), we have

$$\mathcal{F}_0(\phi, p; \phi, p) \geq \left(1 + \left(1 - \frac{1}{\epsilon}\right)2\eta^2\right)\mathcal{S}_0(w_1, w_2; w_1, w_2) + (1 - \epsilon)\gamma\mathcal{S}_0(z, z; z, z).$$

We now make

$$1 + \left(1 - \frac{1}{\epsilon}\right)2\eta^2 = (1 - \epsilon)\gamma$$

by setting

$$\epsilon = \frac{\sqrt{(1 + 2\eta^2 - \gamma)^2 + 8\eta^2\gamma} - (1 + 2\eta^2 - \gamma)}{2\gamma}.$$

Note that $0 < \epsilon < 1$. This and (3.23) yield

$$\mathcal{F}_0(\phi, p; \phi, p) \geq c(\mathcal{S}_0(w_1, w_2; w_1, w_2) + \mathcal{S}_0(z, z; z, z)) = c\mathcal{S}_0(\phi, p; \phi, p),$$

with

$$c = \frac{(1 + 2\eta^2 + \gamma) - \sqrt{(1 + 2\eta^2 + \gamma)^2 - 4\gamma}}{2} > 0,$$

which is (3.21).

Finally, show that $\alpha\mathcal{S} \leq \widehat{\mathcal{F}}$: We now have

$$(3.33) \quad \mathcal{S}_0(\phi, p; \phi, p) = \|\nabla^*A\nabla\phi\|_{0,\Omega}^2 + \|A^{1/2}\nabla p\|_{0,\Omega,n}^2 \leq \frac{1}{c}\mathcal{F}_0(\phi, p; \phi, p).$$

By (3.33) and the definition of \mathcal{F}_0 , we have

$$\begin{aligned} (3.34) \quad \|A^{1/2}\nabla\phi\|_{0,\Omega,n}^2 &\leq 2\left(\|A^{1/2}\nabla p\|_{0,\Omega,n}^2 + \|A^{1/2}\nabla p - A^{1/2}\nabla\phi\|_{0,\Omega,n}^2\right) \\ &\leq 2\left(\frac{1}{c} + 1\right)\mathcal{F}_0(\phi, p; \phi, p). \end{aligned}$$

The Poincaré–Friedrichs inequality (3.1) and the definition of \mathcal{F}_0 again, together with (3.33) and (3.34), now yield

$$\begin{aligned} \mathcal{S}(A\nabla\phi, p; A\nabla\phi, p) &= \|A^{1/2}\nabla\phi\|_{0,\Omega,n}^2 + \|\nabla^*A\nabla\phi\|_{0,\Omega}^2 + \|p\|_{0,\Omega}^2 + \|A^{1/2}\nabla p\|_{0,\Omega,n}^2 \\ &\leq \left(\frac{3 + 2c + d}{c}\right)\widehat{\mathcal{F}}(A\nabla\phi, p; A\nabla\phi, p) \end{aligned}$$

which yields (3.20) with

$$\alpha_0 = \frac{c}{3 + 2c + d},$$

and the theorem is proved. \square

4. Equivalent forms. Problem (2.1)–(2.3) can be written as a first-order system in many different ways. In this section, we illustrate that such formulations often lead to least-squares functionals that are equivalent to the form $G(\mathbf{u}, p; f)$ defined in (2.15). For example, if $X = \mathbf{a}^* \nabla + cI$, then one might choose the system

$$(4.1) \quad \begin{pmatrix} -I & A\nabla \\ \nabla^* + \mathbf{a}^* A^{-1} & cI \end{pmatrix} \begin{pmatrix} \mathbf{u} \\ p \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ f \end{pmatrix}$$

with $\mathbf{u} \in \mathbf{W}$ and $p \in V$. The functional associated with this system would then be

$$(4.2) \quad G_\ell(\mathbf{u}, p; f) = \|A\nabla p - \mathbf{u}\|_{0,\Omega,n}^2 + \|\nabla^* \mathbf{u} + \mathbf{a}^* A^{-1} \mathbf{u} + cp - f\|_{0,\Omega}^2.$$

System (4.1) can be related to the system implicit in $G(\mathbf{u}, p; f)$ by the simple expression

$$(4.3) \quad \begin{pmatrix} I & 0 \\ -\mathbf{a}^* A^{-1} & I \end{pmatrix} \begin{pmatrix} -I & A\nabla \\ \nabla^* & \mathbf{a}^* \nabla + cI \end{pmatrix} \begin{pmatrix} \mathbf{u} \\ p \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ f \end{pmatrix}.$$

Thus, the ratio of $G(\mathbf{u}, p; f)$ and $G_\ell(\mathbf{u}, p; f)$ is bounded above and below, respectively, by the largest and smallest singular values of the transformation matrix. Specifically, if

$$(4.4) \quad \mathbf{a}^* A^{-2} \mathbf{a} \leq C$$

for almost every $x \in L^2(\Omega)$, then

$$(4.5) \quad \alpha_\ell^- G(\mathbf{u}, p; 0) \leq G_\ell(\mathbf{u}, p; 0) \leq \alpha_\ell^+ G(\mathbf{u}, p; 0),$$

for all $\mathbf{u} \in \mathbf{W}$ and $p \in V$, where

$$(4.6) \quad \alpha_\ell^\pm = \frac{(2 + C) \pm \sqrt{C^2 + 4C}}{2}.$$

Note that $0 < \alpha_\ell^- < \alpha_\ell^+ < \infty$. Since equivalence is transitive, then Theorem 3.1 holds with $G(\mathbf{u}, p; 0)$ replaced by $G_\ell(\mathbf{u}, p; 0)$.

If, on the other hand, we have $X = \nabla^* \mathbf{b} + cI$, then one might consider the system

$$(4.7) \quad \begin{pmatrix} -I & A\nabla + \mathbf{b} \\ \nabla^* & cI \end{pmatrix} \begin{pmatrix} \hat{\mathbf{u}} \\ p \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ f \end{pmatrix}.$$

The functional associated with this system would then be

$$(4.8) \quad G_r(\hat{\mathbf{u}}, p; f) = \|A\nabla p + \mathbf{b}p - \hat{\mathbf{u}}\|_{0,\Omega,n}^2 + \|\nabla^* \hat{\mathbf{u}} + cp - f\|_{0,\Omega}^2.$$

If $(\hat{\mathbf{u}}, p)$ satisfies system (4.7), then $\hat{\mathbf{u}} = A\nabla p + \mathbf{b}p$. Boundary condition (2.3) now implies

$$(4.9) \quad \mathbf{n} \cdot \hat{\mathbf{u}} = (\mathbf{n} \cdot \mathbf{b})p \quad \text{on } \Gamma_N.$$

The choice of spaces on which to pose the minimization of G_r is not immediately clear. However, we can again relate system (4.7) to the system associated with G by writing (4.7) as

$$(4.10) \quad \begin{pmatrix} -I & A\nabla \\ \nabla^* & \nabla^*\mathbf{b} + cI \end{pmatrix} \begin{pmatrix} I & -\mathbf{b} \\ 0 & I \end{pmatrix} \begin{pmatrix} \hat{\mathbf{u}} \\ p \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ f \end{pmatrix}.$$

From (4.10) it is clear that

$$(4.11) \quad G_r(\hat{\mathbf{u}}, p; f) = G(\hat{\mathbf{u}} - \mathbf{b}p, p; f).$$

The two functionals become analogous if we pose the minimization of $G_r(\hat{\mathbf{u}}, p; f)$ over the space

$$(4.12) \quad \mathbf{Z} \equiv \{(\hat{\mathbf{u}}, p) \in H(\text{div}) \times V : \mathbf{n} \cdot (\hat{\mathbf{u}} - \mathbf{b}p) = 0 \text{ on } \Gamma_N\}.$$

The ellipticity of G_r with respect to the $H(\text{div}) \times H^1(\Omega)$ norm can now be established by setting $\mathbf{u} = \hat{\mathbf{u}} - \mathbf{b}p$ and noting that

$$(4.13) \quad \frac{1}{K} \left(\|\mathbf{u}\|_{H(\text{div})}^2 + \|p\|_{1,\Omega}^2 \right) \leq \left(\|\hat{\mathbf{u}}\|_{H(\text{div})}^2 + \|p\|_{1,\Omega}^2 \right) \leq K \left(\|\mathbf{u}\|_{H(\text{div})}^2 + \|p\|_{1,\Omega}^2 \right),$$

where $K = 2 + 4\|\nabla^*\mathbf{b}\|_{0,\Omega}^2 + 4\|\mathbf{b}\|_{0,\Omega,n}^2$. Combining (4.11) and (4.13), we see that Theorem 3.1 again holds with G replaced by G_r evaluated on \mathbf{Z} .

System (4.7) has the unfortunate property that the space \mathbf{Z} is not a tensor product space. A vector $(\hat{\mathbf{u}}, p)$ is admissible only if $\hat{\mathbf{u}}$ and p satisfy the proper relationship at the boundary. On the other hand, if boundary condition (2.3) were replaced by

$$(4.14) \quad \mathbf{n} \cdot A\nabla p + bp = 0 \quad \text{on } \Gamma_N,$$

where b is now a scalar, and if a $\mathbf{b} \in H(\text{div})$ could be determined so that

$$(4.15) \quad \mathbf{n} \cdot \mathbf{b} = b \quad \text{on } \Gamma_N,$$

then one might consider the system

$$(4.16) \quad \begin{pmatrix} -I & A\nabla + \mathbf{b} \\ \nabla^* & X - \nabla^*\mathbf{b} \end{pmatrix} \begin{pmatrix} \hat{\mathbf{u}} \\ p \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ f \end{pmatrix},$$

and its associated functional

$$(4.17) \quad G_m(\hat{\mathbf{u}}, p; f) = \|A\nabla p + \mathbf{b}p - \hat{\mathbf{u}}\|_{0,\Omega,n}^2 + \|\nabla^*\hat{\mathbf{u}} + Xp - \nabla^*(\mathbf{b}p) - f\|_{0,\Omega}^2.$$

The proper subspace is once again $\mathbf{W} \times V$ as in (2.7). A slight modification of the proof of Theorem 3.1 yields that G can be replaced by G_m as well.

5. Finite element approximation. We approximate the minimum of $G(\mathbf{u}, p; f)$ in (2.15) using a Rayleigh–Ritz-type finite element method. Let \mathcal{T}_h be a partition of the domain Ω into finite elements, i.e., $\Omega = \cup_{K \in \mathcal{T}_h} K$ with $h = \max\{\text{diam}(K) : K \in \mathcal{T}_h\}$. Assume that the triangulation \mathcal{T}_h is regular (see [12]). Let V_h and \mathbf{W}_h be finite-dimensional subspaces of V and \mathbf{W} with the following properties:

$$(5.1) \quad \inf_{q_h \in V_h} \|q - q_h\|_{1,\Omega} \leq Ch^k \|q\|_{k+1,\Omega},$$

$$(5.2) \quad \inf_{\mathbf{v}_h \in \mathbf{W}_h} \|\mathbf{v} - \mathbf{v}_h\|_{H(\text{div})} \leq Ch^r \|\mathbf{v}\|_{r+1, \Omega},$$

where $k, r > 0$ are integers, $q \in H^{k+1}(\Omega)$, and $\mathbf{v} \in (H^{r+1}(\Omega))^n$.

A standard choice for such spaces is piecewise polynomials of degree k and r , respectively, i.e.,

$$(5.3) \quad V_h = \{q_h \in C^0(\Omega) : q_h|_K \in P_k(K), \forall K \in \mathcal{T}_h, q_h = 0 \text{ on } \Gamma_D\}$$

$$(5.4)$$

$$\mathbf{W}_h = \{\mathbf{v}_h \in H(\text{div}) : (\mathbf{v}_h)_i|_K \in P_r(K), i = 1, \dots, n, \forall K \in \mathcal{T}_h, \mathbf{v}_h \cdot \mathbf{n} = 0 \text{ on } \Gamma_N\}.$$

Here, $P_s(K)$ is the space of polynomials of degree s on K . Another choice for \mathbf{W}_h is a Raviart–Thomas space (see [27]), where continuity requirements are weaker. In this case, of course, definition (5.4) should be modified since the Raviart–Thomas spaces consist of incomplete polynomials of degree r on each element. Also, we allow the use of standard isoparametric elements.

The finite element approximation to minimizing $G(\mathbf{u}, p; f)$ in (2.15) on $\mathbf{W} \times V$ becomes: find $p_h \in V_h$ and $\mathbf{u}_h \in \mathbf{W}_h$ that satisfy

$$(5.5) \quad G(\mathbf{u}_h, p_h; f) = \min_{\mathbf{v}_h \in \mathbf{W}_h, q_h \in V_h} G(\mathbf{v}_h, q_h; f).$$

By Theorem 3.1 and the fact that $V_h \in V$ and $\mathbf{W}_h \in \mathbf{W}$, we conclude that (5.5) has a unique solution and is equivalent to the weak form: find $p_h \in V_h$ and $\mathbf{u}_h \in \mathbf{W}_h$ such that

$$(5.6) \quad \mathcal{F}(\mathbf{u}_h, p_h; \mathbf{v}_h, q_h) = (f, \nabla^* \mathbf{v}_h + X q_h)$$

for every $q_h \in V_h$ and $\mathbf{v}_h \in \mathbf{W}_h$. Moreover, the error satisfies the orthogonality property

$$(5.7) \quad \mathcal{F}(\mathbf{u} - \mathbf{u}_h, p - p_h; \mathbf{v}_h, q_h) = 0$$

for every $q_h \in V_h$ and $\mathbf{v}_h \in \mathbf{W}_h$.

THEOREM 5.1. *Let $s = \min(k, r)$ and assume that $p \in H^{s+1}(\Omega)$ and $\mathbf{u} \in (H^{s+1}(\Omega))^n$. Then*

$$\|p - p_h\|_{1, \Omega} + \|\mathbf{u} - \mathbf{u}_h\|_{H(\text{div})} \leq Ch^s (\|p\|_{s+1, \Omega} + \|\mathbf{u}\|_{s+1, \Omega}),$$

where the constant C does not depend on h , p , or \mathbf{u} .

Proof. The proof is a simple consequence of the orthogonality property (5.7) and the approximation properties (5.1) and (5.2) of the finite element spaces V_h and \mathbf{W}_h . \square

This error estimate is optimal for $k = r$. In many applications, however, it may be useful to have higher-order approximations to the fluxes or velocities \mathbf{u} , which is why we allow $r > k$. Moreover, numerical experience in [9], [26], and [25] indicates that, for certain values of r and k such that $r \neq k$, the error in \mathbf{u} is $O(h^r)$, while the error in p is $O(h^k)$. It may be possible to develop a theory that bounds the errors in \mathbf{u} and p separately.

6. Condition number. In this section, we bound the condition number of the linear system arising from (5.6). To this end, we additionally assume that the triangulation \mathcal{T}_h satisfies the inverse assumption, i.e., there exists a positive constant θ such that, for all $K \in \mathcal{T}_h$,

$$(6.1) \quad h \leq \theta \operatorname{diam}(K).$$

Under assumption (6.1), many standard finite element spaces satisfy the so-called inverse inequality (see [12]); i.e.,

$$(6.2) \quad |q|_{1,\Omega} \leq C h^{-1} \|q\|_{0,\Omega} \quad \forall q \in V_h,$$

$$(6.3) \quad |\operatorname{div} \mathbf{v}|_{0,\Omega} \leq C h^{-1} \|\mathbf{v}\|_{0,\Omega} \quad \forall \mathbf{v} \in \mathbf{W}_h.$$

Let ϕ_1, \dots, ϕ_N and ψ_1, \dots, ψ_M be bases for V_h and \mathbf{W}_h , respectively. Then, for any $q \in V_h$ and $\mathbf{v} \in \mathbf{W}_h$, we have

$$q = \sum_{i=1}^N \eta_i \phi_i \quad \text{and} \quad \mathbf{v} = \sum_{i=1}^M \xi_i \psi_i.$$

Denote by $|\boldsymbol{\eta}|$ and $|\boldsymbol{\xi}|$ the ℓ_2 -norms of the vectors $\boldsymbol{\eta}$ and $\boldsymbol{\xi}$, respectively. Under assumption (6.1), all well-known finite element spaces with the usual choice of bases have the following property: there exist positive constants α_i and β_i ($i = 1, 2$) such that

$$(6.4) \quad \alpha_1 h^n |\boldsymbol{\eta}| \leq \|q\|_{0,\Omega} \leq \alpha_2 h^n |\boldsymbol{\eta}|$$

and

$$(6.5) \quad \beta_1 h^n |\boldsymbol{\xi}| \leq \|\mathbf{v}\|_{0,\Omega} \leq \beta_2 h^n |\boldsymbol{\xi}|.$$

(Recall that n is the spatial dimension of Ω .)

THEOREM 6.1. *Assume that inequalities (6.2)–(6.5) hold. Then the condition number of the linear system resulting from (5.6) is $O(h^{-2})$.*

Proof. It follows from ellipticity and continuity of the bilinear form \mathcal{F} and the inverse inequalities (6.2) and (6.3) that

$$\alpha (\|\mathbf{v}\|_{0,\Omega}^2 + \|q\|_{0,\Omega}^2) \leq \mathcal{F}(\mathbf{v}, q; \mathbf{v}, q) \leq C h^{-2} (\|\mathbf{v}\|_{0,\Omega}^2 + \|q\|_{0,\Omega}^2).$$

Combining inequalities (6.4) and (6.5), we obtain

$$C h^n (|\boldsymbol{\xi}|^2 + |\boldsymbol{\eta}|^2) \leq \mathcal{F}(\mathbf{v}, q; \mathbf{v}, q) \leq C h^{n-2} (|\boldsymbol{\xi}|^2 + |\boldsymbol{\eta}|^2).$$

This completes the proof. \square

REFERENCES

- [1] A. K. AZIZ, R. B. KELLOGG, AND A. B. STEPHENS, *Least-squares methods for elliptic systems*, Math. Comp., 44 (1985), pp. 53–70.
- [2] I. BABUŠKA, *The finite element method with Lagrange multipliers*, Numer. Math., 20 (1973), pp. 179–192.
- [3] P. B. BOCHEV AND M. D. GUNZBURGER, *Accuracy of least-squares methods for the Navier–Stokes equations*, Comput. Fluids, 22 (1993), pp. 549–563.

- [4] P. B. BOCHEV AND M. D. GUNZBURGER, *Analysis of least-squares finite element methods for the Stokes equations*, Math. Comp., to appear.
- [5] F. BREZZI, *On existence, uniqueness and approximation of saddle-point problems arising from Lagrange multipliers*, RAIRO Anal. Numér., 2 (1974), pp. 129–151.
- [6] F. BREZZI, J. DOUGLAS, M. FORTIN, AND D. MATINI, *Efficient rectangular mixed finite elements in two and three space variables*, RAIRO Modél. Math. Anal. Numér., 21 (1987), pp. 581–604.
- [7] F. BREZZI, J. DOUGLAS, AND D. MARINI, *Two families of mixed finite elements for second order elliptic problems*, Numer. Math., 47 (1985), pp. 217–235.
- [8] Z. CAI, T. MANTEUFFEL, AND S. MCCORMICK, *First-order system least-squares for second-order partial differential equations: Part II*, SIAM J. Numer. Anal., 33 (1996), to appear.
- [9] G. F. CAREY AND Y. SHEN, *Convergence studies of least-squares finite elements for first order systems*, Comm. Appl. Numer. Meth., 5 (1989), pp. 427–434.
- [10] T. F. CHEN, *On the least-squares approximations to compressible flow problems*, Numer. Meth. PDE's, 2 (1986), pp. 207–228.
- [11] T. F. CHEN AND G. J. FIX, *Least-squares finite element simulation of transonic flows*, Appl. Numer. Math., 2 (1986), pp. 399–408.
- [12] P. G. CIARLET, *The Finite Element Methods for Elliptic Problems*, North-Holland, Amsterdam, 1978.
- [13] E. D. EASON, *A review of least-squares methods for solving partial differential equations*, Internat. J. Numer. Math. Engrg., 10 (1976), pp. 1021–1046.
- [14] L. FRANCA AND R. STENBERG, *Error analysis of some Galerkin-least-squares method for the elasticity equations*, Report #1054, 1989, INRIA, Le Chesnay, France, pp. 1–21.
- [15] V. GIRAULT AND P. A. RAVIART, *Finite Element Methods for Navier-Stokes Equations: Theory and Algorithms*, Springer-Verlag, New York, 1986.
- [16] C. I. GOLDSTEIN, T. A. MANTEUFFEL, AND S. V. PARTER, *Preconditioning and boundary conditions without H_2 estimates: L_2 condition numbers and the distribution of the singular values*, SIAM J. of Numer. Anal., 13 (1992), pp. 259–288.
- [17] P. GRISVARD, *Elliptic Problems in Nonsmooth Domains*, Pitman, Boston, 1985.
- [18] T. J. R. HUGHES AND L. P. FRANCA, *A new finite element formulation for computational fluid dynamics. VII. The Stokes problems with various well-posed boundary conditions: symmetric formulation that converges for all velocity pressure spaces*, Comput. Meth. Appl. Mech. Engrg., 65 (1987), pp. 85–96.
- [19] T. J. R. HUGHES, L. P. FRANCA, AND M. BULESTRA, *A new finite element formulation for computational fluid dynamics. V. Circumventing the Babuška-Brezzi condition: A stable Petrov-Galerkin formulation of the Stokes problem accomodating equal-order interpolations*, Comput. Meth. Appl. Mech. Engrg., 59 (1986), pp. 85–99.
- [20] D. C. JESPERSEN, *A least-square decomposition method for solving elliptic systems*, Math. Comp., 31 (1977), pp. 873–880.
- [21] B. N. JIANG AND C. CHANG, *Least-squares finite elements for the Stokes problem*, Comput. Meth. Appl. Mech. Engrg., 81 (1990), pp. 13–37.
- [22] B. N. JIANG AND L. A. POVINELLI, *Optimal least-squares finite element method for elliptic problems*, Comput. Meth. Appl. Mech. Engrg., 102 (1993), pp. 199–212.
- [23] O. A. LADYZHENSKAYA, *The Mathematical Theory of Viscous Incompressible Flows*, Gordon and Breach, London, 1969.
- [24] P. NEITTAANMÄKI AND J. SARANEN, *On finite element approximation of the gradient for the solution to Poisson equation*, Numer. Math., 37 (1981), pp. 131–148.
- [25] A. I. PEHLIVANOV, G. F. CAREY, AND R. D. LAZAROV, *Least squares mixed finite elements for second order elliptic problems*, SIAM J. Numer. Anal., 31 (1994), to appear.
- [26] A. I. PEHLIVANOV, G. F. CAREY, R. D. LAZAROV, AND Y. SHEN, *Convergence of least squares finite elements for first order ODE systems*, Computing, 51 (1993), pp. 111–123.
- [27] P. A. RAVIART AND I. M. THOMAS, *A mixed finite element method for second order elliptic problems*, Lecture Notes in Math. 606, Springer-Verlag, Berlin, New York, 1977, pp. 292–315.
- [28] W. L. WENDLAND, *Elliptic Systems in the Plane*, Pitman, London, 1979.