# A Petrov-Galerkin Mixed Finite Element Method with Exponential Fitting

**R. R. P. van Nooyen***
*CWI, Centre for Mathematics and Information Sciences,
P. O. Box 4079, 1009 AB Amsterdam, Holland*

We discuss a Petrov-Galerkin mixed finite element formulation of the semiconductor continuity equations on a rectangular domain. We give error estimates for equations that are in principle degenerate in the singularly perturbed case. We give arguments that indicate that the method is also effective in the singularly perturbed case. We develop a discretization that gives a higher-order accurate solution for use in an *a posteriori* error estimator. © 1995 John Wiley & Sons, Inc.

## I. INTRODUCTION

The use of a form of exponential fitting for the semiconductor continuity equation is suggested by the success of the Scharfetter–Gummel discretization [1] in one dimension and variations on that discretization in two dimensions. Numerous derivations of Scharfetter–Gummel type discretizations are given in the literature, for instance by Selberherr [2], Markowich [3], Bank et al. [4], Brezzi et al. [5], and others. This article extends a one-dimensional exponential fitting technique, discussed by Hemker [6], to the two-dimensional context.

In Section II we introduce a model equation for the semiconductor continuity equations and several bilinear forms, related to the coefficients in this equation. In Sections III and IV we treat the discretization. In Section V we collect some technical results, and in Section VI we derive two error estimates. These error estimates are based on the techniques used by Douglas and Roberts [7]. The proofs in Section VI take all characteristics of our special discrete system into account, in particular the quadrature rule for the approximation of certain integrals in the discrete system. Note that the error estimates in Section VI are degenerate if the problem is singularly perturbed, i.e., if the convection dominates in the problem. On the other hand, an indication for good behavior of the method for singular problems is that—for constant coefficients—it reproduces the solution $C \exp(-\beta_1 x_1 - \beta_2 x_2)$ exactly. In Section IX, we develop an *a posteriori* error estimator, and in the last section we discuss our findings.

*Present address: Fac TWI/Gebouw ET, Ye Etage, Vakgroep Toegepaste Analyse, Delft University of Technology, Mekelweg 4, Delft, Holland.

## II. EQUATION

We consider the following problem, find $u \in H^2(\Omega)$ such that

$$-\mathrm{div}\left(\frac{1}{\alpha}\left(\mathbf{grad}\ u + u\boldsymbol{\beta}\right)\right) + \gamma u = f \ in\ \Omega \quad \text{and}$$

$$u = -g \quad \text{on}\ \partial\Omega, \tag{1}$$

where $\Omega$ is a bounded rectangular domain in $\mathbf{R}^2$. We impose the following restrictions on the coefficients:

$$\alpha \in W_1^\infty(\Omega) \quad \text{and} \quad \exists\ A \in \mathbf{R}\colon \alpha \geq A > 0 \ in\ \Omega, \tag{2}$$

$$\frac{1}{\alpha} \in W_1^\infty(\Omega), \tag{3}$$

$$\boldsymbol{\beta} = (\beta_1, \beta_2)^T \quad \text{with}\ \beta_1, \beta_2 \in W_1^\infty(\Omega), \tag{4}$$

$$\gamma \in W_1^\infty(\Omega) \quad \text{and} \quad \gamma \geq 0 \ in\ \Omega, \tag{5}$$

where $W_1^\infty(\Omega)$, $H^2(\Omega)$ are the usual Sobolev spaces [8], and

$$H(\mathrm{div}, \Omega) := \{\boldsymbol{\tau} \in L^2(\Omega)^2 \mid \mathrm{div}\ \boldsymbol{\tau} \in L^2(\Omega)\},$$

with scalar product

$$(\boldsymbol{\sigma}, \boldsymbol{\tau})_{H(\mathrm{div},\Omega)} = \int_\Omega \boldsymbol{\sigma} \cdot \boldsymbol{\tau}\, d\mu + \int_\Omega \mathrm{div}\ \boldsymbol{\sigma}\ \mathrm{div}\ \boldsymbol{\tau}\, d\mu,$$

is a Hilbert space (see also Girault and Raviart [9] formula 2.15 in Section 2.2). We assume that the equation has a solution and that $f \in L^2(\Omega)$, $g \in H^{3/2}(\partial\Omega)$.

The stationary semiconductor continuity equations take the form (1). Here $\boldsymbol{\beta}$ corresponds to the electric field, the term $\gamma u$ corresponds to a linear approximation to the recombination term, and $1/\alpha$ corresponds to the electron or hole mobility. The exact correspondence depends on the choice of scaling [10].

In order to formulate the weak mixed form of this equation, we use the following bilinear forms

$$(s, t) = \int_\Omega st\, d\mu \quad \forall\ s, t \in L^2(\Omega),$$

$$a(\boldsymbol{\sigma}, \boldsymbol{\tau}) = \int_\Omega \alpha\boldsymbol{\sigma} \cdot \boldsymbol{\tau}\, d\mu \quad \forall\ \boldsymbol{\sigma}, \boldsymbol{\tau} \in H(\mathrm{div}, \Omega),$$

$$b(\boldsymbol{\sigma}, t) = \int_\Omega \boldsymbol{\beta} \cdot \boldsymbol{\sigma}t\, d\mu \quad \forall\ \boldsymbol{\sigma} \in H(\mathrm{div}, \Omega), \quad t \in L^2(\Omega),$$

$$c(s, t) = \int_\Omega \gamma st\, d\mu \quad \forall\ s, t \in L^2(\Omega),$$

$$\langle g, h \rangle = \int_{\partial\Omega} gh\, d\lambda \quad \forall\ g, h \in L^2(\partial\Omega).$$

Given these definitions, we see immediately that any solution $u \in H^2(\Omega)$ of (1) generates a solution $(\boldsymbol{\sigma}, u) \in H(\mathrm{div}, \Omega) \times L^2(\Omega)$ of

$$a(\boldsymbol{\sigma}, \boldsymbol{\tau}) - (\mathrm{div}\ \boldsymbol{\tau}, u) + b(\boldsymbol{\tau}, u) = \langle g, \boldsymbol{\tau} \cdot \mathbf{n}_{\partial\Omega} \rangle \quad \forall\ \boldsymbol{\tau} \in H(\mathrm{div}, \Omega), \tag{6a}$$

$$(\mathrm{div}\ \boldsymbol{\sigma}, t) + c(u, t) = (f, t) \quad \forall\ t \in L^2(\Omega), \tag{6b}$$

where

$$\boldsymbol{\sigma} = -\frac{1}{\alpha} \left( \text{grad } u + u\boldsymbol{\beta} \right).$$

In order to simplify the notation, we denote the Cartesian product of a normed linear space $E$ with itself by $\mathbf{E}$ in boldface type, $\mathbf{E} := E \times E$. We define

$$\|(\mu_1, \mu_2)^T\|_{\mathbf{E}} := \left( \sum_{i=1}^{2} \|\mu_i\|_E^2 \right)^{1/2} \quad \forall \ (\mu_1, \mu_2)^T \in \mathbf{E}.$$

## III. PREPARATIONS

We introduce a partition of the domain and define the adjoint problem of (1), which we use in the derivation of one of our error estimates. Next, we introduce several special projections that are needed in the definition of our approximation spaces and in the derivation of the error estimates. Finally, we give an error estimate for the projections.

### A. Partitioning the Domain

We assume that our domain $\Omega$ is rectangular. On $\Omega$, we use Cartesian coordinates, with the unit vectors $\mathbf{e}_1$ and $\mathbf{e}_2$ parallel to the edges of $\Omega$. We define $\tau_i := \boldsymbol{\tau} \cdot \mathbf{e}_i$ for $\boldsymbol{\tau} \in \mathbf{L}^2(\Omega)$ and $x_i := \mathbf{x} \cdot \mathbf{e}_i$ for $\mathbf{x} \in \mathbf{R}^2$. Before we treat our discretization, we define our approximation space. We assume that our partition is the cartesian product of partitions

$$P = \{0 = p_0 < p_1 < \cdots < p_{N_1} = L_1\}, \tag{7}$$

and

$$Q = \{0 = q_0 < q_1 < \cdots < q_{N_2} = L_2\} \tag{8}$$

of the sides of our domain. We define the index set $K$,

$$K = \{(i + 1/2, j + 1/2) \mid i = 0, 1, \ldots, N_1 - 1, j = 0, 1, \ldots, N_2 - 1\},$$

with the obvious index pair for a given cell,

$$\Omega_{i+1/2, j+1/2} = \{\mathbf{x} \mid p_i < x_1 < p_{i+1}, q_j < x_2 < q_{j+1}\}.$$

We define $\mathbf{x}_k$ to be the center of $\Omega_k$ and $\mathbf{h}_k$ to be the diagonal of $\Omega_k$, with the notation $\chi_k$ for the characteristic function of $\Omega_k$. (The characteristic function of a set is the function that is equal to one in all points of the set and zero elsewhere.) The edges of $\Omega_k$ are the sets:

$$\Gamma_{k,i,j} = \{\mathbf{x} \in \overline{\Omega}_k \mid \mathbf{x} \cdot \mathbf{e}_i = (\mathbf{x}_k + (j - 1/2)\mathbf{h}_k) \cdot \mathbf{e}_i\} \quad \text{for } i = 1, 2, j = 0, 1. \tag{9}$$

$\chi_{k,i,j}$ is the characteristic function of edge $\Gamma_{k,i,j}$. So $(i, j) = (1, 0), (1, 1), (2, 0), (2, 1)$ denote the left, right, bottom and top edges.

## B. Adjoint Problem

We use the following definition for the adjoint problem of (1) (cf. Douglas and Roberts [7]), $w \in H^2(\Omega)$,

$$-\text{div}\left(\frac{1}{\alpha} \text{ grad } w\right) + \frac{\beta}{\alpha} \cdot \text{grad } w + \gamma w = f \text{ in } \Omega, \qquad w = 0 \quad \text{on } \partial\Omega. \qquad (10)$$

The adjoint problem is called regular, if there is a unique solution $w$ for every $f \in L^2(\Omega)$ and this solution satisfies $\|w\|_{H^2(\Omega)} \leq C\|f\|_{L^2(\Omega)}$ for every $f \in L^2(\Omega)$.

Throughout this article the upper case $C$, without a subscript, denotes a generic constant, which may have a different value at each appearance.

The weak mixed form of the adjoint problem is:

$$(\tau, w) \in H(\text{div}, \Omega) \times L^2(\Omega), \qquad (11)$$

$$a(\tau, \sigma) - (\text{div } \sigma, w) = 0 \quad \forall \sigma \in H(\text{div}, \Omega) \quad \text{and} \qquad (11a)$$

$$(\text{div } \tau, t) - b(\tau, t) + c(w, t) = (f, t) \quad \forall t \in L^2(\Omega). \qquad (11b)$$

Any solution $w \in H^2(\Omega)$ of (10) generates a solution $(-1/\alpha \text{ grad } w, w)$ of this problem. If (10) is regular, then this solution satisfies $\|w\|_{H^2(\Omega)} + \|\tau\|_{H^1(\Omega)} \leq C\|f\|_{L^2(\Omega)}$.

## C. Projections

We introduce several local projections, and use these to define four global mappings, $P_h$, $\mathbf{P}_h$, $\Pi_h$, and $\tilde{\Pi}_h$ that map function spaces to finite dimensional function spaces. First, we define $P[\Omega_k]$ to be the orthogonal projection from $L^2(\Omega_k)$ to the space of constant functions on $\Omega_k$, and we define $P[\Gamma_{k,i,j}]$ to be the orthogonal projection from $L^2(\Gamma_{k,i,j})$ to the space of constant functions on $\Gamma_{k,i,j}$.

We use $P[\Omega_k]$ to create two global mappings, $P_h: L^2(\Omega) \rightarrow L^2(\Omega)$,

$$P_h f = \sum_{k \in K} \chi_k P[\Omega_k](f) \quad \forall f \in L^2(\Omega), \qquad (12a)$$

and $\mathbf{P}_h: \mathbf{L}^2(\Omega) \rightarrow \mathbf{L}^2(\Omega)$,

$$\mathbf{P_h}\beta = \sum_{k \in K} \chi_k(P[\Omega_k](\beta \cdot \mathbf{e}_1)\mathbf{e}_1 + P[\Omega_k](\beta \cdot \mathbf{e}_2)\mathbf{e}_2) \quad \forall \beta \in \mathbf{L}^2(\Omega). \qquad (12b)$$

Next, we introduce two mappings, based on $P[\Gamma_{k,i,j}]$. These have as their domain the space $\Sigma$,

$$\Sigma := \{\tau \in H(\text{div}, \Omega) \,|\, \tau|_{\partial\Omega_k} \cdot \mathbf{n}_{\partial\Omega_k} \in L^2(\partial\Omega_k) \quad \forall k \in K\}.$$

This space is similar to that introduced by Roberts and Thomas in formula (1.10) of their report [11].

To simplify the definition of these mappings, we introduce local coordinates on each cell $\Omega_k$,

$$\vec{\xi}_k := \begin{pmatrix} \dfrac{x_1 - x_{k,1}}{h_{k,1}} + \dfrac{1}{2} \\[3mm] \dfrac{x_2 - x_{k,2}}{h_{k,2}} + \dfrac{1}{2} \end{pmatrix}. \qquad (13)$$

The mappings are defined as follows:

$$\Pi_h \tau = \sum_{k \in K} \chi_k \sum_{i=1}^{2} ((1 - \xi_{k,i}) P[\Gamma_{k,i,0}](\tau_i) + \xi_{k,i} P[\Gamma_{k,i,1}](\tau_i)) e_i, \qquad (14)$$

$$\tilde{\Pi}_h \tau = \sum_{k \in K} \chi_k \sum_{i=1}^{2} ((1 - \zeta_{k,i}) P[\Gamma_{k,i,0}](\tau_i) + \zeta_{k,i} P[\Gamma_{k,i,1}](\tau_i)) e_i, \qquad (15)$$

where

$$\zeta_{k,i} = \begin{cases} \dfrac{\exp(\xi_{k,i} h_{k,i} P[\Omega_k](\beta_i)) - 1}{\exp(h_{k,i} P[\Omega_k](\beta_i)) - 1} & \text{if } P[\Omega_k](\beta_i) \neq 0, \\[2mm] \xi_{k,i} & \text{if } P[\Omega_k](\beta_i) = 0. \end{cases}$$

For $\Pi_h \tau$ we get the $i^{th}$ component on $\Omega_k$ by linear interpolation between the projections of this component on the two sides orthogonal to $e_i$. For $\tilde{\Pi}_h \tau$, however, we obtain the same component by using an exponential function to interpolate between the projections of this component on the two sides orthogonal to $e_i$.

The following finite dimensional function spaces are now introduced as the ranges of the above projections:

$$V_h = \Pi_h(\Sigma), \qquad W_h = P_h(L^2(\Omega)), \qquad \text{and} \quad X_h = \tilde{\Pi}_h(\Sigma).$$

$V_h \times W_h$ is the lowest order Raviart–Thomas–Nedelec space for rectangles. This space and the above projections were described by Douglas and Roberts [7], Raviart and Thomas [12], and, for $\Omega \subset \mathbf{R}^3$, by Nedelec [13]. In this article, we use the usual space, $V_h \times W_h$ as the trial function space and $X_h \times W_h$ as the test function space. In effect, we use exponential test functions instead of the usual linear test functions. Thus, we obtain a Petrov–Galerkin mixed finite element discretization.

## D. Error Estimates for Projections

A lemma on the accuracy of our projections is now found. Considering the number and diversity of articles on error estimates, e.g., the classical projection estimates from Ciarlet and Raviart [14], this may seem superfluous, but we shall see that the relative simplicity of the case under consideration makes it possible to derive sharp error estimates under minimal assumptions.

**Lemma 1.** *If $f$ is a square integrable function with square integrable derivatives on a rectangle $\Omega = [0, h_1] \times [0, h_2]$ with sides $\Gamma_{1,1} = \{h_1\} \times [0, h_2]$, $\Gamma_{2,1} = [0, h_1] \times \{h_2\}$, $\Gamma_{1,0} = \{0\} \times [0, h_2]$, and $\Gamma_{2,0} = [0, h_1] \times \{0\}$, then the following inequalities hold:*

$$\| f - P[\Omega]f \|_{L^2(\Omega)} \leq (2h_1^2 + 2h_2^2)^{1/2} \| \mathbf{grad}\, f \|_{L^2(\Omega)}. \qquad (16a)$$

*If $s$ is a continuous function with domain $[0, h_1]$ and range $[0, 1]$, then we have*

$$\| f - (1 - s)\Pi[\Gamma_{1,0}]f - s\Pi[\Gamma_{1,1}]f \|_{L^2(\Omega)} \leq (2h_1^2 + 2h_2^2)^{1/2} \| \mathbf{grad}\, f \|_{L^2(\Omega)}. \qquad (16b)$$

*If $f \in L^\infty(\Omega)$, $\mathbf{grad}\, f \in \mathbf{L}^\infty(\Omega)$, then*

$$\| f - P[\Omega]f \|_{L^\infty(\Omega)} \leq (h_1 + h_2) \| \mathbf{grad}\, f \|_{\mathbf{L}^\infty(\Omega)}. \qquad (16c)$$

**Proof.** We start by proving the above inequalities for $f \in C^1(\Omega)$. We can then extend them by the usual density argument to $H^1(\Omega)$. To prove the first inequality, we write

$$\|f - P[\Omega]f\|^2_{L^2(\Omega)} = \int_{x=0}^{h_1} \int_{y=0}^{h_2} \left( \frac{1}{h_1 h_2} \int_{w=0}^{h_1} \int_{z=0}^{h_2} f(x,y) - f(w,z) \, dw \, dz \right)^2 dx \, dy,$$

by definition,

$$f(x,y) - f(w,z) = \int_{a=w}^{x} \frac{\partial f}{\partial a}(a,z) \, da + \int_{b=z}^{y} \frac{\partial f}{\partial b}(x,b) \, da.$$

If we substitute this into the above expression, then we find

$$\|f - P[\Omega]f\|^2_{L^2(\Omega)} =$$

$$\int_{x=0}^{h_1} \int_{y=0}^{h_2} \left( \frac{1}{h_1 h_2} \int_{w=0}^{h_1} \int_{z=0}^{h_2} \left[ \int_{a=w}^{x} \frac{\partial f}{\partial a}(a,z) \, da + \int_{b=z}^{y} \frac{\partial f}{\partial b}(x,b) \, db \right] dw \, dz \right)^2 dx \, dy.$$

We apply the Hölder inequality to the inner integrals and extend the integrations over $a$ and $b$, where appropriate,

$$\|f - P[\Omega]f\|^2_{L^2(\Omega)} \leq$$

$$\int_{x=0}^{h_1} \int_{y=0}^{h_2} \left( \frac{h_1^{1/2}}{h_2^{1/2}} \|\partial f/\partial x_1\|_{L^2(\Omega)} + h_2^{1/2} \left( \int_{b=0}^{h_2} \left( \frac{\partial f}{\partial b}(x,b) \right)^2 db \right)^{1/2} \right)^2 dx \, dy.$$

We use $(|A| + |B|)^2 \leq 2(A^2 + B^2)$ to write this as

$$\|f - P[\Omega]f\|^2_{L^2(\Omega)} \leq 2 \int_{x=0}^{h_1} \int_{y=0}^{h_2} \frac{h_1}{h_2} \|\partial f/\partial x_1\|^2_{L^2(\Omega)} \, dx \, dy$$

$$+ 2 \int_{y=0}^{h_2} h_2 \|\partial f/\partial x_2\|^2_{L^2(\Omega)} \, dy.$$

This reduces to,

$$\|f - P[\Omega]f\|^2_{L^2(\Omega)} \leq 2h_1^2 \|\partial f/\partial x_1\|^2_{L^2(\Omega)} + 2h_2^2 \|\partial f/\partial x_2\|^2_{L^2(\Omega)}.$$

Now, we consider the second inequality, (16b), we write,

$$\|f - (1 - s)\Pi[\Gamma_{1,0}]f - s\Pi[\Gamma_{1,1}]f\|^2_{L^2(\Omega)} =$$

$$\int_{x=0}^{h_1} \int_{y=0}^{h_2} \left( \frac{1}{h_2} \int_{z=0}^{h_2} [(1 - s(x))(f(x,y) - f(0,z)) + s(x)(f(x,y) - f(h_1,z))] dz \right)^2 dx \, dy.$$

We use partial derivatives to rewrite the expression,

$$\|f - (1 - s)\Pi[\Gamma_{1,0}]f - s\Pi[\Gamma_{1,1}]f\|^2_{L^2(\Omega)} =$$

$$\int_{x=0}^{h_1} \int_{y=0}^{h_2} \left( \frac{1}{h_2} \int_{z=0}^{h_2} \left[ (1 - s(x)) \left\{ \int_{a=0}^{x} \frac{\partial f}{\partial a}(a,z) \, da + \int_{b=z}^{y} \frac{\partial f}{\partial b}(x,b) \, db \right\} \right. \right.$$

$$\left. \left. + s(x) \left\{ \int_{a=h_1}^{x} \frac{\partial f}{\partial a}(a,z) \, da + \int_{b=z}^{y} \frac{\partial f}{\partial b}(x,b) \, db \right\} \right] dz \right)^2 dx \, dy.$$

Using the Hölder inequality we extend the integrals where appropriate, so that

$$\| f - (1 - s)\Pi[\Gamma_{1,0}]f - s\Pi[\Gamma_{1,1}]f \|^2_{L^2(\Omega)} \le$$

$$\int_{x=0}^{h_1} \int_{y=0}^{h_2} \left( \frac{h_1^{1/2}}{h_2^{1/2}} \| \partial f/\partial x_1 \|_{L^2(\Omega)} + h_2^{1/2} \left( \int_{b=0}^{h_2} \left( \frac{\partial f}{\partial b}(x,b) \right)^2 db \right)^{1/2} \right)^2 dx\, dy\, .$$

We use $(|A| + |B|)^2 \le 2(A^2 + B^2)$ to write this as

$$\| f - (1 - s)\Pi[\Gamma_{1,0}]f - s\Pi[\Gamma_{1,1}]f \|^2_{L^2(\Omega)} \le$$

$$2 \int_{x=0}^{h_1} \int_{y=0}^{h_2} \left( \frac{h_1}{h_2} \| \partial f/\partial x_1 \|_{L^2(\Omega)} + h_2 \int_{b=0}^{h_2} \left( \frac{\partial f}{\partial b}(x,b) \right)^2 db \right) dx\, dy\, ,$$

which reduces to

$$\| f - (1 - s)\Pi[\Gamma_{1,0}]f - s\Pi[\Gamma_{1,1}]f \|^2_{L^2(\Omega)} \le 2h_1^2 \| \partial f/\partial x_1 \|^2_{L^2(\Omega)} + 2h_2^2 \| \partial f/\partial x_2 \|^2_{L^2(\Omega)}\, .$$

Lastly we verify (16c),

$$f(x,y) - f(w,z) = \int_{a=w}^{x} \frac{\partial f}{\partial a}(a,z)\, da + \int_{b=z}^{y} \frac{\partial f}{\partial b}(x,b)\, da\, .$$

So,

$$\frac{1}{h_1 h_2} \int_{x=0}^{h_1} \int_{y=0}^{h_2} f(x,y) - f(w,z)\, dx\, dy =$$

$$\frac{1}{h_1 h_2} \int_{x=0}^{h_1} \int_{y=0}^{h_2} \left( \int_{a=w}^{x} \frac{\partial f}{\partial a}(a,z)\, da + \int_{b=z}^{y} \frac{\partial f}{\partial b}(x,b)\, da \right) dx\, dy$$

$$\le (h_1 + h_2)\, \| \mathbf{grad}\, f \|_{L^\infty(\Omega)}\, .$$

■

Note that the above inequalities imply

$$\| \boldsymbol{\sigma} - \Pi_h \boldsymbol{\sigma} \|_{L^2(\Omega)} \le \max_{k \in K}(2h_{k,1}^2 + 2h_{k,2}^2)^{1/2} \| \boldsymbol{\sigma} \|_{H^1(\Omega)}\, , \tag{17a}$$

$$\| \boldsymbol{\sigma} - \tilde{\Pi}_h \boldsymbol{\sigma} \|_{L^2(\Omega)} \le \max_{k \in K}(2h_{k,1}^2 + 2h_{k,2}^2)^{1/2} \| \boldsymbol{\sigma} \|_{H^1(\Omega)}\, , \tag{17b}$$

$$\| u - P_h u \|_{L^2(\Omega)} \le \max_{k \in K}(2h_{k,1}^2 + 2h_{k,2}^2)^{1/2} \| u \|_{H^1(\Omega)}\, , \tag{17c}$$

for suitable $u$ and $\boldsymbol{\sigma}$.

## IV. DISCRETIZATION

We describe our discretization. The basic idea of mixed finite elements with a lowest order Raviart–Thomas trial space and exponentially fitted test subspace for the vector valued functions is complicated by the use of a quadrature rule, needed to keep the M-matrix property for the system without Lagrange multipliers for nonzero $\gamma$. This quadrature rule is discussed in Section IVA. Another complication is the approximation of the coefficients by piecewise constant functions, as described below. In Section IVB we give the resulting discretization.

We replace the coefficients $\alpha$, $\beta$ and $\gamma$ by two-dimensional step functions. To write our modified problem in weak form, we need to define three new bilinear forms:

$$\bar{a}(\sigma, \tau) = \int_\Omega \sigma \cdot \tau P_h \alpha \, d\mu \quad \forall \, \sigma, \tau \in \Sigma,$$

$$\bar{b}(\sigma, t) := \int_\Omega t\sigma \cdot P_h \beta \, d\mu \quad \forall \, \sigma \in \Sigma, t \in L^2(\Omega),$$

$$\bar{c}(s, t) := \int_\Omega st P_h \gamma \, d\mu \quad \forall \, s, t \in L^2(\Omega).$$

The bar on the bilinear forms denotes that the coefficients are replaced by their cell-wise averages. We then replace $\bar{a}$ by $\bar{a}_q$, the subscript $q$ indicates that a—not yet specified—quadrature rule will be used in the evaluation of this bilinear form.

## A. Quadrature Rule

We construct a quadrature rule $\bar{a}_{h,1}$ by imposing the condition that, if $\alpha, \beta$ are constant, $\gamma \equiv 0$, and the solution satisfies $u = C \exp(-\beta_1 x_1 - \beta_2 x_2) + K$, with $C, K \in \mathbf{R}$, then the discrete solution should satisfy $\sigma_h = \Pi_h \sigma$ and $u_h = P_h u$. We see that for the $u$ given above $\sigma = -K\beta/\alpha$, so $\sigma$ is constant. We define $\alpha_h$ separately for each basis function $\eta_{i,j+1/2}$ where

$$\eta_{i,j+1/2} = \begin{cases} \zeta_{(i-1/2,j+1/2)}\mathbf{e}_1 \in \Omega_{i-1/2,j+1/2}, \\ (1 - \zeta_{(i+1/2,j+1/2)})\mathbf{e}_1 \in \Omega_{i+1/2,j+1/2}, \\ 0 \text{ elsewhere}, \end{cases}$$

and $\eta_{i+1/2,j}$, where

$$\eta_{i+1/2,j} = \begin{cases} \zeta_{(i+1/2,j-1/2)}\mathbf{e}_2 \in \Omega_{i-1/2,j-1/2}, \\ (1 - \zeta_{(i+1/2,j-1/2)})\mathbf{e}_2 \in \Omega_{i+1/2,j+1/2}, \\ 0 \text{ elsewhere}. \end{cases}$$

We denote the set of all possible indices for the basis functions $\eta$ by

$$E = \{e = (i, j - 1/2) \mid i = 0, 1, 2, \ldots, N_1, j = 1, 2, \ldots, N_2\} \bigcup$$

$$\{e = (i - 1/2, j) \mid i = 1, 2, \ldots, N_1, j = 0, 1, 2, \ldots, N_2\}.$$

Our quadrature rule should satisfy the following condition:

$$\bar{a}_{h,1}(\sigma, \eta_r) = \bar{a}(\sigma, \eta_r), \tag{A}$$

for all constant $\sigma$ and all $r \in E$. Due to our assumption that the coefficients are constant, we have $a = \bar{a}$ and $b = \bar{b}$. The above condition guarantees that for constant coefficients and constant $\sigma$,

$$a(\sigma, \tau_h) - (u, \text{div } \tau_h) + b(\tau_h, u) = \bar{a}_{h,1}(\Pi_h \sigma, \tau_h) - (P_h u, \text{div } \tau_h)$$

$$+ b(\tau_h, P_h u) \quad \forall \, \tau_h \in X_h,$$

and we also have

$$(\text{div } \sigma, t) = (\text{div } \Pi_h \sigma, t) = 0 \quad \forall \, t \in L^2(\Omega).$$

So our condition (A) on $\overline{a}_{h,1}$ is sufficient for our purposes. We now select the quadrature rule by taking the following definition for $\overline{a}_{h,1}$,

$$\overline{a}_{h,1}(\sigma, \tau) = \sum_{k \in K} \sum_{i=1}^{2} \mu(\Omega_k) P[\Omega_k](\alpha) \left( P[\Omega_k](\zeta_{k,i}) P[\Gamma_{k,i,1}](\sigma_i \tau_i) \right.$$

$$\left. + P[\Omega_k](1 - \zeta_{k,i}) P[\Gamma_{k,i,0}](\sigma_i \tau_i) \right). \quad (18a)$$

We introduce a new problem dependent norm on $X_h$

$$\|\tau_h\|_h = \sum_{k \in K} \sum_{i=1}^{2} \mu(\Omega_k) \left( P[\Omega_k](\zeta_{k,i}) P[\Gamma_{k,i,1}](\tau_{h,i}^2) \right.$$

$$\left. + P[\Omega_k](1 - \zeta_{k,i}) P[\Gamma_{k,i,0}](\tau_{h,i}^2) \right)^{1/2}. \quad (18b)$$

From this point onward, we take $\overline{a}_q = \overline{a}_{h,1}$.

## B. Discrete System

We approximate the solution $(\sigma, u)$ of (6) by $(\sigma_h, u_h) \in V_h \times W_h$, where

$$\overline{a}_q(\sigma_h, \tau) - (u_h, \text{div } \tau) + \overline{b}(\tau, u_h) = \langle \tau \cdot \eta_{\partial\Omega}, g \rangle \quad \forall \tau \in X_h, \quad (19a)$$

$$(\text{div } \sigma_h, t) + \overline{c}(u_h, t) = (f, t) \quad \forall t \in W_h. \quad (19b)$$

If we use $\overline{a}$ instead of $\overline{a}_q$, then this means that our discrete problem does not always yield an M-matrix for $u_h$. Consider, for instance, the corresponding discretization on a uniform mesh with mesh width $h$ in one dimension with $\alpha = 1$, $\beta = \vec{0}$, and $\gamma$ constant. If $\alpha \gamma h^2 / 6 > 1$, then the off-diagonal elements of the discretization matrix for $u_h$ after elimination of $\sigma_h$ through static condensation have the same sign as the elements on the diagonal.

The idea of using linear trial functions and exponential test functions was used by Hemker for singularly perturbed two-point boundary problems [6]. For the one-dimensional case, the introduction of exponential test functions follows from the requirement that it must be possible to approximate the Green's function of the problem by the test functions.

In the following sections, we prove that the solution of our discretization (19) is an $O(h)$ approximation to the solution of our original problem.

## V. TECHNICAL RESULTS

This section contains some technical results, collected for later reference.

**Lemma 2.**

$$\tilde{\Pi}_h \circ \Pi_h = \tilde{\Pi}_h, \quad (20a)$$

$$\Pi_h \circ \tilde{\Pi}_h = \Pi_h, \quad (20b)$$

$$(\text{div } \sigma, P_h t) = (\text{div } \Pi_h \sigma, t) \quad \forall \sigma \in \Sigma, t \in L^2(\Omega), \quad (20c)$$

$$\Pi_h \tau \cdot \vec{n}_{\partial\Omega} = \tilde{\Pi}_h \tau \cdot \vec{n}_{\partial\Omega} \quad \forall \tau \in \Sigma. \quad (20d)$$

**Proof.** Both mappings are based on the same projections $P[\Gamma_{k,i,j}]$, so (20a) and (20b) are trivial.

To prove (20c) we use a special case of Green's theorem:

$$\int_{\Omega_k} \operatorname{div} \sigma \, d\mu = \sum_{i=1}^{2} \frac{\mu(\Omega_k)}{h_{k,i}} \left( P[\Gamma_{k,i,1}](\sigma_i) - P[\Gamma_{k,i,0}](\sigma_i) \right).$$

If we combine this with the definition of $\Pi_h$, the proof of (20c) is complete. Equation (20d) follows immediately from the definitions.  ■

**Lemma 3.** *If* $\sigma \in \Sigma$ *and we define* $a_{k,i} = P[\Gamma_{k,i,0}](\sigma_i)$ *and* $b_{k,i} = P[\Gamma_{k,i,1}](\sigma_i)$, *then the following inequalities hold for* $\|\Pi_h \sigma\|_{L^2(\Omega_k)}$ *and* $\|\tilde{\Pi}_h \sigma\|_{L^2(\Omega_k)}$:

$$\frac{\mu(\Omega_k)}{6} \sum_{i=1}^{2} (a_{k,i}^2 + b_{k,i}^2) \leq \|\Pi_h \sigma\|_{L^2(\Omega_k)}^2 \leq \frac{\mu(\Omega_k)}{2} \sum_{i=1}^{2} (a_{k,i}^2 + b_{k,i}^2). \tag{21a}$$

$$\|\tilde{\Pi}_h \sigma\|_{L^2(\Omega_k)}^2 \leq 2\|\tilde{\Pi}_h \sigma\|_h^2 \leq 12\|\Pi_h \sigma\|_{L^2(\Omega_k)}^2. \tag{21b}$$

**Proof.** Formula (21a) follows immediately from

$$(\Pi_h \sigma, \Pi_h \sigma) = \sum_{i=1}^{2} \sum_{k \in K} \int_{\Omega_k} ((1 - \xi_{k,i})a_{k,i} + \xi_{k,i} b_{k,i})^2 \, d\mu.$$

Next, we derive (21b) from

$$(\tilde{\Pi}_h \sigma, \tilde{\Pi}_h \sigma) = \sum_{i=1}^{2} \sum_{k \in K} \int_{\Omega_k} ((1 - \zeta_{k,i})a_{k,i} + \zeta_{k,i} b_{k,i})^2 \, d\mu.$$

We see immediately that

$$\int_{\Omega_k} ((1 - \zeta_{k,i})a_{k,i} + \zeta_{k,i} b_{k,i})^2 \, d\mu \leq \int_{\Omega_k} 2(1 - \zeta_{k,i})^2 a_{k,i}^2 + 2\zeta_{k,i}^2 b_{k,i}^2 \, d\mu \leq$$

$$2\int_{\Omega_k} (1 - \zeta_{k,i})a_{k,i}^2 + \zeta_{k,i} b_{k,i}^2 \, d\mu = 2\mu(\Omega_k)(P[\Omega_k](1 - \zeta_{k,i})a_{k,i}^2 + P[\Omega_k](\zeta_{k,i})b_{k,i}^2).$$

This implies (21b).  ■

Lemma 4 shows, that $\bar{a}$ is $L^2(\Omega)$-bounded and $L^2(\Omega)$-elliptic.

**Lemma 4.** *Let* $\alpha \in W_1^\infty(\Omega)$, $\alpha \geq A > 0 \in \Omega$ *and* $\bar{a}(\sigma, \tau) := \int_\Omega P_h(\alpha)\sigma \cdot \tau \, d\mu$ $\forall \sigma, \tau \in L^2(\Omega)$, *then*

$$\bar{a}(\sigma, \tau) \leq \|\alpha\|_{L^\infty(\Omega)} \|\sigma\|_{L^2(\Omega)} \|\tau\|_{L^2(\Omega)} \quad \forall \sigma, \tau \in L^2(\Omega), \tag{22a}$$

*and*

$$\bar{a}(\tau, \tau) \geq A\|\tau\|_{L^2(\Omega)} \quad \forall \tau \in L^2(\Omega). \tag{22b}$$

**Proof.** From (2) it follows that

$$A \leq \frac{\int_{\Omega_k} \alpha \, d\mu}{\mu(\Omega_k)} \leq \|\alpha\|_{L^\infty(\Omega)},$$

together with the definitions of $P$ and $\bar{a}$, this implies (22a) and (22b).  ■

We introduce the minimum mesh width $h_{min}$ and the maximum mesh width $h_{max}$,

$$h_{min} = \min_{k \in K} \min_{i=1,2} |h_{k,i}|, \tag{23a}$$

$$h_{max} = \max_{k \in K} \max_{i=1,2} |h_{k,i}|. \tag{23b}$$

## A. Properties of $\bar{a}_q$

We discuss the properties of the quadrature rule $\bar{a}_q$ and assume that $\bar{a}_q = \bar{a}_{h,1}$, where $\bar{a}_{h,1}$ is given by (18a). We also discuss the interaction between $\Pi$, $\tilde{\Pi}$, and $\bar{a}_q$. We show that $\bar{a}_q$ is $L^2(\Omega)$-bounded on $V_h$, and we also show that $\bar{a}_q$ is $L^2(\Omega)$-elliptic on $V_h$ and $X_h$.

**Lemma 5.** *If* $\sigma, \tau \in \Sigma$, *then*

$$\bar{a}_q(\Pi_h\sigma, \Pi_h\tau) = \bar{a}_q(\Pi_h\tau, \Pi_h\sigma) = \bar{a}_q(\sigma, \Pi_h\tau) = \bar{a}_q(\Pi_h\sigma, \tau) = \bar{a}_q(\sigma, \tilde{\Pi}_h\tau) =$$

$$\bar{a}_q(\tilde{\Pi}_h\sigma, \tau) = \bar{a}_q(\tilde{\Pi}_h\sigma, \tilde{\Pi}_h\tau), \quad (24a)$$

$$\|\alpha\|_{L^\infty(\Omega)}\|\tilde{\Pi}_h\sigma\|_h^2 \geq \bar{a}_q(\tilde{\Pi}_h\sigma, \tilde{\Pi}_h\sigma) \geq \frac{1}{2}\bar{a}(\tilde{\Pi}_h\sigma, \tilde{\Pi}_h\sigma) \geq \frac{A}{2}\|\tilde{\Pi}_h\sigma\|_{L^2(\Omega)}^2, \quad (24b)$$

$$\bar{a}_q(\Pi_h\sigma, \Pi_h\tau) \leq 6\|\alpha\|_{L^\infty(\Omega)}\|\Pi_h\sigma\|_{L^2(\Omega)}\|\tilde{\Pi}_h\tau\|_h, \quad (24c)$$

$$A\|\tilde{\Pi}_h\tau\|_h^2 \leq \bar{a}_q(\tau, \tilde{\Pi}_h\tau) \leq \|\alpha\|_{L^\infty(\Omega)}\|\tilde{\Pi}_h\tau\|_h^2. \quad (24d)$$

**Proof.** The definitions of $\Pi_h$, $\tilde{\Pi}_h$, and $\bar{a}_q$ imply (24a). Inequality (24b) follows immediately from (18a), (18b), and (21b). To prove (24c), we need some auxiliary variables, $a_{k,i} = P[\Gamma_{k,i,0}](\sigma)$, $b_{k,i} = P[\Gamma_{k,i,1}](\sigma)$, $c_{k,i} = P[\Gamma_{k,i,0}](\tau)$, and $d_{k,i} = P[\Gamma_{k,i,1}](\tau)$. Cauchy–Schwarz is used twice, we obtain

$$\bar{a}_q(\Pi\sigma, \tilde{\Pi}\tau) = \sum_{k \in K} P[\Omega_k](\alpha)\mu(\Omega_k) \sum_{i=1}^{2} (P[\Omega_k](1 - \zeta_{k,i})a_{k,i}c_{k,i} + P[\Omega_k](\zeta_{k,i})b_{k,i}d_{k,i})$$

$$\leq \sum_{k \in K} P[\Omega_k](\alpha)\mu(\Omega_k)\left(\sum_{i=1}^{2}(a_{k,i}^2 + b_{k,i}^2)\right)^{1/2}$$

$$\times \left(\sum_{i=1}^{2}[P[\Omega_k](1 - \zeta_{k,i})^2 c_{k,i}^2 + P[\Omega_k](\zeta_{k,i})^2 d_{k,i}^2)\right)^{1/2}.$$

We use

$$P[\Omega_k](f)^2 \leq P[\Omega_k](f^2)$$

to rewrite the term in $c$ and $d$, and we use (21a) to replace the term in $a$ and $b$ by $\|\Pi_h\sigma\|_{L^2(\Omega)}$,

$$\bar{a}_q(\Pi_h\sigma, \tilde{\Pi}_h\tau) \leq \sum_{k \in K} P[\Omega_k](\alpha)\mu(\Omega_k)6 \frac{\|\Pi_h\sigma\|_{L^2(\Omega)}}{\mu(\Omega_k)^{1/2}}$$

$$\times \left(\sum_{i=1}^{2}(P[\Omega_k]((1 - \zeta_{k,i})^2)c_{k,i}^2 + P[\Omega_k]((\zeta_{k,i})^2)d_{k,i}^2)\right)^{1/2}.$$

We see immediately that this implies

$$\bar{a}_q(\Pi_h\sigma, \tilde{\Pi}_h\tau) \leq 6\|\alpha\|_{L^\infty(\Omega)}\|\Pi_h\sigma\|_{L^2(\Omega)}\|\tilde{\Pi}_h\tau\|_h.$$

This proves (24c). Inequality (24d) follows immediately from (18). ∎

## B. Difference between $\bar{a}$ and $\bar{a}_q$

For our error estimates, we need an upper bound for the difference between the value of $a(\sigma_h, \tau)$ and that of $\bar{a}_q(\sigma_h, \tau)$ for $\sigma_h \in V_h$, $\tau \in \mathbf{H}^1(\Omega)$. As we already know from (16c) (see also Lemmas 8 and 9) that

$$|a(\sigma, \tau_h) - \bar{a}(\sigma, \tau_h)| \leq 2h_{\max} \|\alpha\|_{W_1^\infty(\Omega)} \|\sigma\|_{L^2(\Omega)} \|\tau_h\|_{L^2(\Omega)},$$

an estimate for $|\bar{a}(\sigma, \tau_h) - \bar{a}_q(\sigma, \tau_h)|$ suffices. Such an estimate is derived in Lemma 6.

**Lemma 6.** *Let* $\tau_h \in X_h$ *and* $\sigma \in \mathbf{H}^1(\Omega)$, *then*

$$|\bar{a}(\sigma, \tau_h) - \bar{a}_q(\sigma, \tau_h)| \leq 2\|\alpha\|_{L^\infty(\Omega)} h_{\max} \|\tau_h\|_h \|\sigma\|_{H^1(\Omega)}. \tag{25}$$

**Proof.** To simplify our notation, we introduce $a_{k,i} = P[\Gamma_{k,i,0}](\tau_h)$, $b_{k,i} = P[\Gamma_{k,i,1}](\tau_h)$, $\sigma_{k,i,0} = P[\Gamma_{k,i,0}](\sigma_i)$, and $\sigma_{k,i,1} = P[\Gamma_{k,i,1}](\sigma_i)$. We prove the lemma for $\sigma$ with $\sigma_1, \sigma_2 \in C^1(\Omega)$, and extend by density.

We consider the difference between the two forms on one subdomain $\Omega_k$ with $P[\Omega_k](\alpha) = 1$.

$$\left| \int_{\Omega_k} \sigma \cdot \tau_h \, d\mu - \mu(\Omega_k) \sum_{i=1}^2 (P[\Omega_k](1 - \zeta_{k,i})P[\Gamma_{k,i,0}](\sigma_i \tau_{h,i}) \right.$$

$$\left. + P[\Omega_k](\zeta_{k,i})P[\Gamma_{k,i,1}](\sigma_i \tau_{h,i})) \right| = \left| \int_{\Omega_k} \sum_{i=1}^2 ((1 - \zeta_{k,i})a_{k,i} + \zeta_{k,i}b_{k,i})\sigma_i \, d\mu - \mu(\Omega_k) \right.$$

$$\times \left. \sum_{i=1}^2 (P[\Omega_k](1 - \zeta_{k,i})P[\Gamma_{k,i,0}](a_{k,i}\sigma_i) + P[\Omega_k](\zeta_{k,i})P[\Gamma_{k,i,1}](b_{k,i}\sigma_i)) \right| =$$

$$\left| \int_{\Omega_k} \sum_{i=1}^2 ((1 - \zeta_{k,i})a_{k,i}\sigma_i + \zeta_{k,i}b_{k,i}\sigma_i - P[\Omega_k](1 - \zeta_{k,i})a_{k,i}\sigma_{k,i,0} \right.$$

$$\left. - P[\Omega_k](\zeta_{k,i})b_{k,i}\sigma_{k,i,1}) \, d\mu \right| = \left| \int_{\Omega_k} \sum_{i=1}^2 ((1 - \zeta_{k,i})a_{k,i}(\sigma_i - \sigma_{k,i,0}) \right.$$

$$\left. + \zeta_{k,i}b_{k,i}(\sigma_i - \sigma_{k,i,1})) \, d\mu \right|.$$

The application of the Cauchy–Schwarz inequality to this last term and insertion of $\alpha$ yields the following result:

$$|\bar{a}(\sigma, \tau_h) - \bar{a}_q(\sigma, \tau_h)| \leq h_{\max} \|\alpha\|_{L^\infty(\Omega)} \|\tau_h\|_h \left( \sum_{k \in K} \sum_{i=1}^2 \sum_{j=0}^1 \|\sigma_i - \sigma_{i,k,j}\|_{L^2(\Omega_k)}^2 \right)^{1/2}.$$

If we take $s = j$ in (16b) then this implies

$$|\bar{a}(\sigma, \tau_h) - \bar{a}_q(\sigma, \tau_h)| \leq \|\alpha\|_{L^\infty(\Omega)} \|\tau_h\|_h \left( \sum_{i=1}^2 4h_{\max}^2 \|\mathrm{grad}\ \sigma_i\|_{L^2(\Omega)}^2 \right)^{1/2}$$

$$\leq 2h_{\max} \|\alpha\|_{L^\infty(\Omega)} \|\tau_h\|_h \|\sigma\|_{H^1(\Omega)}.$$

Because $C^1(\overline{\Omega})$ is dense in $\mathbf{H}^1(\Omega)$, the formula also holds for $\sigma_1, \sigma_2 \in \mathbf{H}^1(\Omega)$. ∎

## VI. ERROR ESTIMATES

We use the standard estimates for $\|\sigma - \Pi_h\sigma\|_{L^2(\Omega)}$ and $\|u - P_h u\|_{L^2(\Omega)}$, as described in Section IIID, to reduce the problem to deriving bounds for $\|P_h u - u_h\|_{L^2(\Omega)}$ and $\|\Pi_h\sigma - \sigma_h\|_{L^2(\Omega)}$. We discuss two possible derivations of an $O(h)$ error bound. The first needs the assumption that $h_{\max}$ is "small enough," while the second places a condition on an approximation of the discrete version of the adjoint problem.

### A. Errors Due to Approximation of the Bilinear Forms

As preparation for the derivation of *a priori* error estimates, we derive some upper bounds on the errors caused by the piecewise constant approximation of the coefficients $\alpha$, $\beta$, and $\gamma$. We use the following well-known notation. If $V$ and $W$ are normed linear spaces, then $L(V, W; \mathbf{R})$ is the space of bounded bilinear forms on $V$ and $W$, the standard norm of an element $b \in L(V, W; \mathbf{R})$ is given by

$$\|b\|_{L(V, W; \mathbf{R})} = \sup_{v \in V} \sup_{w \in W} \frac{|b(v, w)|}{\|v\|_V \|w\|_W} .$$

**Lemma 7.** *If* $\alpha \in W_1^\infty(\Omega)$, *then*

$$\|a - \bar{a}_q\|_{L(H^1(\Omega),(X_h\|\cdot\|_h); \mathbf{R})} \leq 6h_{\max}\|\alpha\|_{W_1^\infty(\Omega)} ,$$

*where* $(X_h, \|\cdot\|_h)$ *is a normed linear space with, as elements, the elements of* $X_h$ *but with* $\|\cdot\|_h$ *as norm.*

**Proof.** From Eqs. (16c) and (21b) it follows that

$$|a(\sigma, \tau_h) - \bar{a}(\sigma, \tau_h)| \leq 4h_{\max}\|\alpha\|_{W_1^\infty(\Omega)}\|\sigma\|_{L^2(\Omega)}\|\tau_h\|_h .$$

When combined with Lemma 6, this implies

$$\|a - \bar{a}_q\|_{L(H^1(\Omega),(X_h,\|\cdot\|_h); \mathbf{R})} \leq 6h_{\max}\|\alpha\|_{W_1^\infty(\Omega)} .$$

∎

**Lemma 8.** *If* $\beta \in W_1^\infty(\Omega)$, *then*

$$\|b - \bar{b}_l\|_{L(L^2(\Omega), L^2(\Omega); \mathbf{R})} \leq 4h_{\max}\|\beta\|_{W_1^\infty(\Omega)} .$$

**Proof.** This follows immediately from (16c).    ∎

**Lemma 9.** *If* $\gamma \in W_1^\infty(\Omega)$, *then*

$$\|c - \bar{c}_q\|_{L(L^2(\Omega), L^2(\Omega); \mathbf{R})} \leq 2h_{\max}\|\gamma\|_{W_1^\infty(\Omega)} .$$

**Proof.** This follows immediately from (16c).    ∎

### B. An *a priori* Error Estimate

The following two lemmas show nice properties of our discretization, which are needed to derive the error bound.

**Lemma 10.**  *Let* $\tau \in \Sigma$, $t \in L^2(\Omega)$, *then*

$$\overline{b}(\tilde{\Pi}_h \tau, t - P_h t) - (\text{div } \tilde{\Pi}_h \tau, t - P_h t) = 0. \tag{26}$$

**Proof.**  A straightforward calculation shows that $\mathbf{P}_h(\boldsymbol{\beta}) \cdot \tilde{\Pi}_h \tau - \text{div } \tilde{\Pi}_h \tau$ is constant on $\Omega_k$. From this, (26) easily follows.    ∎

**Lemma 11.**  *If* $(\sigma, u)$ *is a solution of* (6) *and* $(\sigma_h, u_h)$ *is a solution of* (19), *then*

$$(\text{div}(\sigma - \sigma_h), P_h t) + c(u - u_h, P_h t) = 0 \quad \forall\, t \in L^2(\Omega). \tag{27}$$

**Proof.**  We take (19b),

$$(\text{div } \sigma_h, P_h t) + \overline{c}(u_h, P_h t) = (f, P_h t),$$

$\overline{c}$ is derived by orthogonal $L^2(\Omega_k)$ projection, so this implies

$$(\text{div } \sigma_h, P_h t) + c(u_h, P_h t) = (f, P_h t).$$

If we subtract this from (6b), $(\text{div } \sigma, P_h t) + c(u, P_h t) = (f, P_h t)$, then we find (27).    ∎

We are now ready to give an estimate for $\|\Pi_h \sigma - \sigma_h\|_h$.

**Theorem 1.**  *If* $(\sigma, u)$ *is the solution of* (6), $(\sigma_h, u_h)$ *is the solution of* (19) *and* $(\sigma, u) \in \mathbf{H}^1(\Omega) \times H^2(\Omega)$, *then there exist positive real numbers* $C$ *and* $D$ *such that*

$$C \le \frac{12}{A} \max(1, \|\alpha\|_{W_1^\infty(\Omega)}, \|\beta\|_{W_1^\infty(\Omega)}, \|\gamma\|_{W_1^\infty(\Omega)}) \max(1, \|\sigma\|_{H^1(\Omega)}, \|u\|_{L^2(\Omega)}), \tag{28}$$

$$D \le 2 \frac{\|\beta\|_{L^\infty(\Omega)}}{A},$$

$$\|\Pi_h \sigma - \sigma_h\|_h^2 \le Ch_{\max}(\|\Pi_h \sigma - \sigma_h\|_h + \|P_h u - u_h\|_{L^2(\Omega)})$$

$$+ D\|\Pi_h \sigma - \sigma_h\|_h \|P_h u - u_h\|_{L^2(\Omega)}.$$

**Proof.**  According to (24d), $A\|\tilde{\Pi}_h(\sigma - \sigma_h)\|_h^2 \le \overline{a}_q(\sigma - \sigma_h, \tilde{\Pi}_h(\sigma - \sigma_h))$. This is the starting point for the derivation of our error bound. Equations (6a) and (19a) imply that

$$\overline{a}_q(\sigma - \sigma_h, \tilde{\Pi}_h(\sigma - \sigma_h)) = (\overline{a}_q - a)(\sigma, \tilde{\Pi}_h(\sigma - \sigma_h)) + a(\sigma, \tilde{\Pi}_h(\sigma - \sigma_h))$$

$$- \overline{a}_q(\sigma_h, \tilde{\Pi}_h(\sigma - \sigma_h))$$

$$= (\overline{a}_q - a)(\sigma, \tilde{\Pi}_h(\sigma - \sigma_h)) + (\text{div } \tilde{\Pi}_h(\sigma - \sigma_h), u)$$

$$- b(\tilde{\Pi}_h(\sigma - \sigma_h), u) + \langle g, \mathbf{n}_{\partial\Omega} \cdot \tilde{\Pi}_h(\sigma - \sigma_h)\rangle$$

$$+ \overline{b}(\tilde{\Pi}_h(\sigma - \sigma_h), u_h) - (\text{div } \tilde{\Pi}_h(\sigma - \sigma_h), u_h)$$

$$- \langle g, \mathbf{n}_{\partial\Omega} \cdot \tilde{\Pi}_h(\sigma - \sigma_h)\rangle$$

$$= (\overline{a}_q - a)(\sigma, \tilde{\Pi}_h(\sigma - \sigma_h)) + (\text{div } \tilde{\Pi}_h(\sigma - \sigma_h), u)_{L^2(\Omega)}$$

$$- (b - \overline{b})(\tilde{\Pi}_h(\sigma - \sigma_h), u) - \overline{b}(\tilde{\Pi}_h(\sigma - \sigma_h), u)$$

$$+ \overline{b}(\tilde{\Pi}_h(\sigma - \sigma_h), u_h) - (\text{div } \tilde{\Pi}_h(\sigma - \sigma_h), u_h),$$

where we give $b - \overline{b}$, $\overline{a}_q - a$, etc., their obvious meaning. If we use Lemma 10, we find

$$A\|\tilde{\Pi}_h(\sigma - \sigma_h)\|_h^2 \le (\overline{a}_q - a)(\sigma, \tilde{\Pi}_h(\sigma - \sigma_h)) - (b - \overline{b})(\tilde{\Pi}_h(\sigma - \sigma_h), u)$$

$$+ (\text{div } \tilde{\Pi}_h(\sigma - \sigma_h), P_h u - u_h)_{L^2(\Omega)} - \overline{b}(\tilde{\Pi}_h(\sigma - \sigma_h), P_h u - u_h).$$

If we use (20b) and (20c) to prepare the way, then the application of Lemma 11 to this expression results in

$$A\|\tilde{\Pi}_h(\sigma - \sigma_h)\|_h^2 \le (\bar{a}_q - a)(\sigma, \tilde{\Pi}_h(\sigma - \sigma_h)) - (b - \bar{b})(\tilde{\Pi}_h(\sigma - \sigma_h), u)$$
$$- c(u - u_h, P_h u - u_h) - \bar{b}(\tilde{\Pi}_h(\sigma - \sigma_h), P_h u - u_h).$$

As $\gamma$ is nonnegative according to (5), we may add $c(P_h u - u_h, P_h u - u_h)$ on both sides of the inequality, we find

$$A\|\tilde{\Pi}_h(\sigma - \sigma_h)\|_h^2 + c(P_h u - u_h, P_h u - u_h) \le (\bar{a}_q - a)(\sigma, \tilde{\Pi}_h(\sigma - \sigma_h))$$
$$-(b - \bar{b})(\tilde{\Pi}_h(\sigma - \sigma_h), u) - (c - \bar{c})(u - P_h u, P_h u - u_h)$$
$$- \bar{b}(\tilde{\Pi}_h(\sigma - \sigma_h), P_h u - u_h).$$

We use Lemmas 7, 8, and 9 to reduce this to

$$A\|\tilde{\Pi}_h(\sigma - \sigma_h)\|_h^2 \le h_{\max}(6\|\alpha\|_{W_I^x(\Omega)}\|\sigma\|_{H^1(\Omega)} + 4\|\beta\|_{W_I^x(\Omega)}\|u\|_{L^2(\Omega)})\|\tilde{\Pi}_h(\sigma - \sigma_h)\|_h$$
$$+ 2h_{\max}\|\gamma\|_{W_I^x(\Omega)}\|u - P_h u\|_{L^2(\Omega)}\|P_h u - u_h\|_{L^2(\Omega)}$$
$$+ 2\|\beta\|_{L^x(\Omega)}\|\tilde{\Pi}_h(\sigma - \sigma_h)\|_h\|P_h u - u_h\|_{L^2(\Omega)}.$$

Note that for all $u \in L^2(\Omega)$, $\|u - P_h u\|_{L^2(\Omega)} \le \|u\|_{L^2(\Omega)}$, and $\|\Pi_h \sigma\|_h = \|\tilde{\Pi}_h \sigma\|_h$.  ∎

Next, we prepare for the second part of our error estimate.

**Lemma 12.**  *If* $(\sigma, u)$ *is the solution of* (6), $(\sigma_h, u_h)$ *is a solution of* (19), *and* $(\tau, q)$ *is the solution of the adjoint problem for an arbitrary right-hand side* $p \in L^2(\Omega)$, *then*

$$(\operatorname{div} \tau, P_h u - u_h) - b(\tau, P_h u - u_h) = a(\sigma, \tilde{\Pi}_h \tau) - \bar{a}_q(\sigma_h, \tilde{\Pi}_h \tau)$$
$$+ (b - \bar{b})(\tilde{\Pi}_h \tau, u) + \bar{b}(\tilde{\Pi}_h \tau - \tau, P_h u - u_h) + (\bar{b} - b)(\tau, P_h u - u_h).$$

**Proof.**  We start by replacing $b$ by $\bar{b}$,

$$(\operatorname{div} \tau, P_h u - u_h) - b(\tau, P_h u - u_h) = (\operatorname{div} \tau, P_h u - u_h) - \bar{b}(\tau, P_h u - u_h)$$
$$+ (\bar{b} - b)(\tau, P_h u - u_h).$$

We use (20a) and (20c) to get

$$(\operatorname{div} \tau, P_h u - u_h) - b(\tau, P_h u - u_h) = (\operatorname{div} \tilde{\Pi}_h \tau, P_h u - u_h) - \bar{b}(\tilde{\Pi}_h \tau, P_h u - u_h)$$
$$+ \bar{b}(\tilde{\Pi}_h \tau - \tau, P_h u - u_h) + (\bar{b} - b)(\tau, P_h u - u_h),$$

Lemma 10 to find

$$(\operatorname{div} \tau, P_h u - u_h) - b(\tau, P_h u - u_h) = (\operatorname{div} \tilde{\Pi}_h \tau, u - u_h) - \bar{b}(\tilde{\Pi}_h \tau, u - u_h)$$
$$+ \bar{b}(\tilde{\Pi}_h \tau - \tau, P_h u - u_h) + (\bar{b} - b)(\tau, P_h u - u_h),$$

and Eqs. (6a) and (19a) to produce

$$(\operatorname{div} \tau, P_h u - u_h) - b(\tau, P_h u - u_h) = a(\sigma, \tilde{\Pi}_h \tau) - \langle g, \tilde{\Pi}_h \tau \cdot \mathbf{n}_{\partial\Omega} \rangle - \bar{a}_q(\sigma_h, \tilde{\Pi}_h \tau)$$
$$+ \langle g, \tilde{\Pi}_h \tau \cdot \mathbf{n}_{\partial\Omega} \rangle + (b - \bar{b})(\tilde{\Pi}_h \tau, u) + \bar{b}(\tilde{\Pi}_h \tau - \tau, P_h u - u_h)$$
$$+ (\bar{b} - b)(\tau, P_h u - u_h).$$

∎

**Lemma 13.** *If $(\sigma, u)$ is the solution of* (6), $(\sigma_h, u_h)$ *is a solution of* (19), *and* $(\tau, w)$ *is the solution of the adjoint problem for an arbitrary right-hand side* $p \in L^2(\Omega)$, *then*

$$c(P_h w, u - u_h) = -a(\tau, \tilde{\Pi}_h(\sigma - \sigma_h)) + \bar{b}(\tilde{\Pi}_h(\sigma - \sigma_h), w - P_h w).$$

**Proof.** From Lemma 11,

$$c(P_h w, u - u_h) = -(\text{div}(\sigma - \sigma_h), P_h w),$$

and according to (20b) and (20c) we can rewrite the right-hand side,

$$c(P_h w, u - u_h) = -(\text{div } \tilde{\Pi}_h(\sigma - \sigma_h), P_h w).$$

We wish to use Eq. (26) from Lemma 10 to remove $P_h$. To do this, we must add and subtract a term $\bar{b}(\tilde{\Pi}_h(\sigma - \sigma_h), P_h w)$ on the right-hand side of our equation. We apply Lemma 10 and gather terms in $\bar{b}$ together,

$$c(P_h w, u - u_h) = -(\text{div } \tilde{\Pi}_h(\sigma - \sigma_h), w) + \bar{b}(\tilde{\Pi}_h(\sigma - \sigma_h), w - P_h w).$$

Finally, we use (11a),

$$c(P_h w, u - u_h) = -a(\tau, \tilde{\Pi}_h(\sigma - \sigma_h)) + \bar{b}(\tilde{\Pi}_h(\sigma - \sigma_h), w - P_h w).$$

∎

**Theorem 2.** *Assume the adjoint problem* (11) *has a unique solution for all square integrable right-hand sides, and assume that there is a constant $C_r$ such that, if $(\tau, w)$ is the solution of* (11) *for a given right-hand side $f$, then*

$$\|\tau\|_{H^1(\Omega)} + \|w\|_{H^1(\Omega)} \le C_r \|f\|_{L^2(\Omega)}.$$

*Now, if $(\sigma, u) \in H^1(\Omega) \times H^2(\Omega)$ is the solution of* (6), *and $(\sigma_h, u_h)$ is a solution of* (19), *then there are constants*

$$0 < C, D, E \le 4C_r(1 + 2h_{\max}) \max(\|\alpha\|_{W_1^\infty(\Omega)}, \|\beta\|_{W_1^\infty(\Omega)}, \|\gamma\|_{W_1^\infty(\Omega)}),$$

*such that*

$$\|P_h u - u_h\|_{L^2(\Omega)} \le Ch_{\max}(\|u\|_{L^2(\Omega)} + \|\sigma\|_{H^1(\Omega)}) + Dh_{\max}\|\tilde{\Pi}_h(\sigma - \sigma_h)\|_h$$
$$+ Eh_{\max}\|P_h u - u_h\|_{L^2(\Omega)}.$$

**Proof.** If we have an estimate for $(P_h u - u_h, p)$ for all $p \in L^2(\Omega)$, then we can use

$$\|t\|_{L^2(\Omega)} = \sup_{p \in L^2(\Omega), p \ne 0} \frac{(p, t)}{\|p\|_{L^2(\Omega)}},$$

to find $\|P_h u - u_h\|_{L^2(\Omega)}$. We use the regularity of the adjoint problem (11) to find a solution $(\tau, w) \in H^1(\Omega) \times L^2(\Omega)$ of (11) for a given right-hand side $p \in L^2(\Omega)$. We may write

$$(p, P_h u - u_h) = (\text{div } \tau, P_h u - u_h) - b(\tau, P_h u - u_h) + c(w, P_h u - u_h).$$

If we apply Lemma 12, we find

$$(p, P_h u - u_h) = a(\sigma, \tilde{\Pi}_h \tau) - \bar{a}_q(\sigma_h, \tilde{\Pi}_h \tau) + (b - \bar{b})(\tilde{\Pi}_h \tau, u)$$
$$+ \bar{b}(\tilde{\Pi}_h \tau - \tau, P_h u - u_h) + (\bar{b} - b)(\tau, P_h u - u_h)$$
$$+ c(w - P_h w, P_h u - u_h) + c(P_h w, P_h u - u_h).$$

Using Lemma 13, we find that

$$(p, P_h u - u_h) = a(\sigma, \tilde{\Pi}_h \tau) - \bar{a}_q(\sigma_h, \tilde{\Pi}_h \tau) + (b - \bar{b})(\tilde{\Pi}_h \tau, u)$$
$$+ \bar{b}(\tilde{\Pi}_h \tau - \tau, P_h u - u_h) + (\bar{b} - b)(\tau, P_h u - u_h)$$
$$+ c(w - P_h w, P_h u - u_h) - a(\tau, \tilde{\Pi}_h(\sigma - \sigma_h))$$
$$+ \bar{b}(\tilde{\Pi}_h(\sigma - \sigma_h), w - P_h w),$$

which can be written as

$$(p, P_h u - u_h) = (a - \bar{a}_q)(\sigma, \tilde{\Pi}_h \tau) + \bar{a}_q(\sigma - \sigma_h, \tilde{\Pi}_h \tau) + (b - \bar{b})(\tilde{\Pi}_h \tau, u)$$
$$+ \bar{b}(\tilde{\Pi}_h \tau - \tau, P_h u - u_h) + (\bar{b} - b)(\tau, P_h u - u_h)$$
$$+ c(w - P_h w, P_h u - u_h) - a(\tau, \tilde{\Pi}_h(\sigma - \sigma_h))$$
$$+ \bar{b}(\tilde{\Pi}_h(\sigma - \sigma_h), w - P_h w).$$

We use (24a) to write this as

$$(p, P_h u - u_h) = (a - \bar{a}_q)(\sigma, \tilde{\Pi}_h \tau) - (a - \bar{a}_q)(\tau, \tilde{\Pi}_h(\sigma - \sigma_h)) + (b - \bar{b})(\tilde{\Pi}_h \tau, u)$$
$$+ \bar{b}(\tilde{\Pi}_h \tau - \tau, P_h u - u_h) + (\bar{b} - b)(\tau, P_h u - u_h)$$
$$+ c(w - P_h w, P_h u - u_h) + \bar{b}(\tilde{\Pi}_h(\sigma - \sigma_h), w - P_h w).$$

Use of the regularity of the adjoint problem (11), Lemmas 7, 8, and 9, and the projection error estimates (16a,b,c) leads to

$$\|P_h u - u_h\|_{L^2(\Omega)} \leq C_r(1 + 2h_{max})2h_{max}\|\alpha\|_{W_1^\ast(\Omega)}(\|\sigma\|_{H^1(\Omega)} + \|\tilde{\Pi}_h(\sigma - \sigma_h)\|_h)$$
$$+ 4C_r h_{max}\|\beta\|_{W_1^\ast(\Omega)}(1 + 2h_{max})\|u\|_{L^2(\Omega)} + 2C_r h_{max}\|\beta\|_{L^\ast(\Omega)}\|P_h u - u_h\|_{L^2(\Omega)}$$
$$+ 2C_r h_{max}(h_{max}\|\gamma\|_{W_1^\ast(\Omega)}\|P_h u - u_h\|_{L^2(\Omega)} + \|\beta\|_{L^\ast(\Omega)}\|\tilde{\Pi}_h(\sigma - \sigma_h)\|_{L^2(\Omega)}).$$

This can be written as

$$\|P_h u - u_h\|_{L^2(\Omega)} \leq \tilde{C} h_{max}(1 + h_{max}) + \tilde{D} h_{max}(1 + h_{max})\|\tilde{\Pi}_h(\sigma - \sigma_h)\|_h$$
$$+ \tilde{E} h_{max}(1 + h_{max})\|P_h u - u_h\|_{L^2(\Omega)}.$$

∎

If $h_{max}$ is small enough, Theorems 1 and 2 together give an $O(h_{max})$ error estimate.

An important limit on $h_{max}$ is implied by the form of the estimates in Theorems 1 and 2. The main problem is that large values of $\|\alpha\|_{W_1^\ast(\Omega)}$, $\|\beta\|_{W_1^\ast(\Omega)}$, and $\|\gamma\|_{W_1^\ast(\Omega)}$ decrease the range of $h_{max}$ for which the estimate is valid. Specifically, in Theorem 2 we need to bring the term with coefficient $E$ to the left-hand side, so we need, for example,

$$4h_{max}C_r(1 + 2h_{max}) \max(\|\alpha\|_{W_1^\ast(\Omega)}, \|\beta\|_{W_1^\ast(\Omega)}, \|\gamma\|_{W_1^\ast(\Omega)}) \leq 1/2.$$

Note that $C_r$, the condition of the dual problem, may also depend on $\alpha$.

If the above inequality holds, then Theorem 2 implies

$$\|P_h u - u_h\|_{L^2(\Omega)} \leq 2Ch_{max}(\|u\|_{L^2(\Omega)} + \|\sigma\|_{H^1(\Omega)}) + 2Dh_{max}\|\tilde{\Pi}_h \sigma - \sigma_h\|_h,$$

with $C, D$ as in Theorem 2. Insertion of this estimate in Theorem 1 leads to a similar problem. Here we need an assumption of the form

$$16h_{max}C_r(1 + 2h_{max})\max(\|\alpha\|_{W_1^*(\Omega)}, \|\beta\|_{W_1^*(\Omega)}, \|\gamma\|_{W_1^*(\Omega)})\frac{\|\beta\|_{L^*(\Omega)}}{A} \leq 1/2$$

to bring the square of the error from the right-hand side to the left-hand side. Theorem 1 then gives

$$\|\tilde{\Pi}_h \sigma - \sigma_h\|_h^2 \leq 2Kh_{max}(\|\tilde{\Pi}_h \sigma - \sigma_h\|_h + 2Ch_{max}(\|u\|_{L^2(\Omega)} + \|\sigma\|_{H^1(\Omega)})$$
$$+ 2Dh_{max}\|\tilde{\Pi}_h \sigma - \sigma_h\|_h),$$

with $K$ the constant $C$ from Theorem 1 and $C, D$ as in Theorem 2. If $\|\tilde{\Pi}_h \sigma - \sigma_h\|_h$ is smaller than $h_{max}$, then we find

$$\|\tilde{\Pi}_h \sigma - \sigma_h\|_h^2 \leq 2Kh_{max}(h_{max} + 2Ch_{max}(\|u\|_{L^2(\Omega)} + \|\sigma\|_{H^1(\Omega)}) + 2Dh_{max}),$$

i.e., an $O(h_{max})$ error estimate for $\|\tilde{\Pi}_h \sigma - \sigma_h\|_h$ and, if it is larger than $h_{max}$, then we may divide terms in the right-hand side either by $\|\tilde{\Pi}_h \sigma - \sigma_h\|_h$ or by $h_{max}$, so

$$\|\tilde{\Pi}_h \sigma - \sigma_h\|_h \leq 2Kh_{max}(1 + 2C(\|u\|_{L^2(\Omega)} + \|\sigma\|_{H^1(\Omega)}) + 2Dh_{max}),$$

so again we find an $O(h_{max})$ error estimate for $\|\tilde{\Pi}_h \sigma - \sigma_h\|_h$.


## C. Another Approach

To improve our estimate of $\|P_h u - u_h\|_{L^2(\Omega)}$, we consider the adjoint of the discrete problem. This means, that we look for $(\tau_h, v_h) \in X_h \times W_h$, such that

$$\bar{a}_q(\tau_h, \sigma_h) - (\text{div } \sigma_h, v_h) = 0 \quad \forall \, \sigma_h \in V_h, \tag{29a}$$

$$(\text{div } \tau_h, t_h) - \bar{b}(\tau_h, t_h) + \bar{c}(v_h, t_h) = (f, t_h) \quad \forall \, t_h \in W_h. \tag{29b}$$

We call this system regular, if there is at least one solution for each $f \in P_h(L^2(\Omega))$, and that all solutions for a particular $f$ satisfy

$$\|\tau_h\|_h + \|v_h\|_{L^2(\Omega)} \leq C\|P_h f\|_{L^2(\Omega)}, \tag{29c}$$

with $C$ independent of the mesh size. This is a somewhat less stringent regularity condition than that given for the continuous adjoint problem (10). Note, that $\tau_h \in X_h$, so $\tau_{h,i}$ is a piecewise exponential function on $\Omega_k$ for $i = 1, 2$.

An example of a general condition under which this system is regular is the following:

$$\alpha \geq A > 0, \qquad \gamma \geq C_0 > 0, \qquad \text{and} \quad AC_0 - \|\beta\|_{L^*(\Omega)}^2 \geq C_1 > 0. \tag{30}$$

To show this, we need the following relations:

$$\int_\Omega \frac{P_h(\alpha)}{4} \tau_h \cdot \tau_h - P_h(\beta) \cdot \tau_h v_h + P_h(\gamma)v_h v_h \, d\mu = \tag{31}$$

$$\int_\Omega \frac{P_h(\alpha)}{4} \left( \tau_h - \frac{2P_h(\beta)}{P_h(\alpha)} v_h \right)^2 + \left( P_h(\gamma) - \frac{P_h(\beta)^2}{P_h(\alpha)} \right) v_h v_h \, d\mu \geq \qquad (31a)$$

$$\int_\Omega P_h(\gamma) \left( v_h - \frac{P_h(\beta) \cdot \tau_h}{2P_h(\alpha)} \right)^2 + \left( P_h(\alpha) - \frac{P_h(\beta)^2}{P_h(\gamma)} \right) \frac{\tau_h \cdot \tau_h}{4} \, d\mu . \qquad (31b)$$

We know, that $(\text{div } \tilde{\Pi}_h\sigma, P_h t) = (\text{div } \Pi_h\sigma, P_h t)$, so, if we take the sum of (29a) and (29b) with $\sigma = \Pi_h\tau_h$ and $t = v_h$, we find

$$\bar{a}_q(\tau_h, \Pi_h\tau_h) - \bar{b}(\tau_h, v_h) + \bar{c}(v_h, v_h) = (f, v_h) . \qquad (32)$$

According to (24a), $\bar{a}_q(\tau_h, \Pi_h\tau_h) = \bar{a}_q(\tilde{\Pi}_h\tau_h, \tilde{\Pi}_h\tau_h)$, and by (24b) we have

$$\frac{1}{4} \bar{a}(\tilde{\Pi}_h\sigma, \tilde{\Pi}_h\sigma) \leq \bar{a}_q(\tilde{\Pi}_h\sigma, \tilde{\Pi}_h\sigma) .$$

Hence,

$$\int_\Omega \frac{P_h(\alpha)}{4} \tau_h \cdot \tau_h - P_h(\beta) \cdot \tau_h v_h + P_h(\gamma) v_h v_h \, d\mu \leq \int_\Omega P_h(f) v_h \, d\mu . \qquad (33)$$

This expression is identical to (31), so (31a) is smaller than $(f, v_h)$; combined with (30), this implies

$$\frac{C_1}{A} \|v_h\|_{L^2(\Omega)} \leq \|f\|_{L^2(\Omega)} . \qquad (34a)$$

In the same way, we find that (31b) is smaller than $(f, v_h)$; together with (30) and (34a), this implies

$$\frac{C_1}{(AC_0)^{1/2}} \|\tau_h\|_{L^2(\Omega)} \leq \|f\|_{L^2(\Omega)} . \qquad (34b)$$

From (32) we see that this implies

$$A\|\tau_h\|_h^2 \leq a_q(\tau_h, \tau_h) \leq \|f\|_{L^2(\Omega)}\|v_h\|_{L^2(\Omega)} + \|\beta\|_{L^\infty(\Omega)}\|\tau_h\|_{L^2(\Omega)}\|v_h\|_{L^2(\Omega)}$$
$$+ \|\gamma\|_{L^\infty(\Omega)}\|v_h\|_{L^2(\Omega)}^2 ,$$

this implies that there is a $C$ such that

$$\|\tau_h\|_h \leq C\|f\|_{L^2(\Omega)} .$$

**Theorem 3.**  *If we assume that (29c) holds, then*

$$\|P_h u - u_h\|_{L^2(\Omega)} \leq h_{\max}(6\|\alpha\|_{W_1^\infty(\Omega)}\|\sigma\|_{H^1(\Omega)} + 2(\|\beta\|_{W_1^\infty(\Omega)} + \|\gamma\|_{W_1^\infty(\Omega)}\|u\|_{L^2(\Omega)}) .$$
$$(35)$$

**Proof.**  We use (29b),

$$(P_h u - u_h, P_h f) = (\text{div } \tau_h, P_h u - u_h) - \bar{b}(\tau_h, P_h u - u_h) + \bar{c}(P_h u - u_h, v_h) .$$

Hence, according to Lemma 10 and the definition of $\bar{c}$,

$$(P_h u - u_h, P_h f) = (\text{div } \tau_h, u - u_h) - \bar{b}(\tau_h, u - u_h) + \bar{c}(u - u_h, v_h) .$$

We use (6a) and (19a) to find

$$(P_h u - u_h, P_h f) = (\text{div } \tau_h, u - u_h) - (\overline{b} - b)(\tau_h, u) - b(\tau_h, u) + \overline{b}(\tau_h, u_h)$$
$$+ \overline{c}(u - u_h, v_h)$$
$$= (\overline{b} - b)(\tau_h, u) + a(\sigma, \tau_h) - \overline{a}_q(\sigma_h, \tau_h) + (\overline{c} - c)(u, v_h)$$
$$+ c(u - u_h, v_h).$$

According to (24a) and Lemma 11, this implies

$$(P_h u - u_h, P_h f) = (\overline{b} - b)(\tau_h, u) + (a - \overline{a}_q)(\sigma, \tau_h) + \overline{a}_q(\Pi_h \sigma - \sigma_h, \tau_h)$$
$$+ (\overline{c} - c)(u, v_h) - (\text{div}(\Pi_h \sigma - \sigma_h), v_h).$$

Now, (29a) implies

$$(P_h u - u_h, P_h f) = (\overline{b} - b)(\tau_h, u) + (a - \overline{a}_q)(\sigma, \tau_h) + (\overline{c} - c)(u, v_h).$$

Finally, we use Lemmas 7, 8, and 9 and (29c) to obtain our error estimate (35).    ■

## VII. VERIFICATION OF THE LOCAL MAXIMUM PRINCIPLE

We use the discrete adjoint problem to show that, for this quadrature rule, the matrix after elimination of $\sigma$ by static condensation is an M-matrix. The discrete adjoint problem is defined in (29).

We assume a regular uniform mesh. We denote the matrix corresponding to (29), after elimination of $\sigma_h$, by $A$. We see that the matrix $A$ has nonpositive off-diagonal elements. We shall show that $A$ is an M-matrix. To do this, we use Theorem 5.12, Chapter 5, page 124 of [15]. This theorem states that, for irreducible matrices with nonpositive off-diagonal elements, the M-matrix property is equivalent to the existence of a positive vector with a nonnegative image that is not identically zero. In our case, the vector $(1, 1, \ldots, 1)^T$ has such an image, because all row sums are nonnegative, and any row corresponding to an edge or corner has a positive row-sum.

The fact that the matrix $A$ is irreducible follows from Theorem 3.6 [15], which states that, for a square matrix, irreducibility is equivalent to its di-graph being strongly connected. Inspection shows that the di-graph of the matrix under consideration is indeed strongly connected.

According to Theorem 5.6 [15], $A^T$ is an M-matrix, too. This implies that the discrete equations for the original $u_h$ satisfy a local maximum principle.

The M-matrix property implies that the system for $u_h$ has a unique solution. From the form of the equations for $\sigma_h$, we see that a given $u_h$ induces a unique $\sigma_h$. This implies that our system is always uniquely solvable. A quick calculation of the coefficients of $u_h$ in (19a) shows that, for constant coefficients and large $\beta$, i.e., with large convection diffusion ratios, we get a relation between $\sigma_h$ and $u_h$, where the "upwind" point is weighed more heavily. If $\beta/\alpha$ remains bounded and we go to the limit $|\beta_1| + |\beta_2| \rightarrow \infty$, then we get a first-order upwind scheme. This suggests that the scheme, in which the coefficients are continuously dependent on this ratio, remains useful close to such a limit.

## VIII. A POSTERIORI ESTIMATOR

We use a special quadrature rule and obtain a higher-order discretization. We seek an $\bar{a}_{h,3}(\cdot,\cdot)$, that minimizes $\bar{a} - \bar{a}_{h,3}$. To do this, we choose a special quadrature rule for each $\bar{a}(\cdot,\eta)$, where $\eta$ is one of the basis functions introduced earlier. Due to the nature of our test functions, the quadrature rule is essentially a one-dimensional rule.

### A. Derivation of the Rule

For $\eta_{i,j+1/2}$, we proceed as follows. We replace the two-dimensional integral by a repeated integral, we integrate exactly in the $e_2$ direction, and then use a three-point rule to approximate the remaining integral. As nodes for the last integration, we take either the centers of $\Gamma_{i-1/2,j+1/2,0}$, $\Gamma_{i+1/2,j+1/2,0}$, and $\Gamma_{i+1/2,j+1/2,1}$. Or, if we are at a boundary, the edge center on the boundary and the two next closest edge centers. We choose the weights as follows:

$$\alpha_h(\Pi_h\sigma,\eta_{i,j+1/2}) = \alpha(\sigma,\eta_{i,j+1/2}),$$

for all $\sigma$ with $x_1$-components that are second-order polynomials in $x_1$, i.e., for all $a,b,c \in \mathbf{R}$, and all $\eta_{i,j+1/2}$, we have

$$\alpha_h(\Pi_h((ax_1^2 + bx_1 + c)e_1),\eta_{i,j+1/2}) = \alpha((ax_1^2 + bx_1 + c)e_1,\eta_{i,j+1/2}).$$

In a similar manner, we define the rule for $\eta_{i+1/2,j}$.

### B. Estimator for the Local Discretization Error and a Lower Bound for the Global Error

We use this rule to obtain an *a posteriori* estimator for the local discretization error and a lower bound for the global error. It is immediately obvious that

$$\bar{a}_{h,3}(\sigma,\eta_r) - \bar{a}_{h,1}(\sigma,\eta_r) \geq O(h_{\max}^2),$$

where $r$ is a possible index-tuple. Moreover,

$$\bar{a}(\sigma,\eta_r) - \bar{a}_{h,3}(\sigma,\eta_r) = O(h_{\max}^3),$$

if $\sigma$ is smooth enough. If

$$|\bar{a}_{h,3}(\rho_h,\eta_r) - (\text{div } \eta_r - \beta \cdot \eta_r,w_h)| \geq K,$$

then we have either

$$\|w_h\|_{L^\infty(\Omega)} \geq C_1 K,$$

or

$$\|\rho_h\|_{L^\infty(\Omega)} \geq C_2 K.$$

We see immediately that, if $(\sigma_h, u_h)$ is the solution of (19) with $\bar{a}_q = \bar{a}_{h,1}$, then

$$\bar{a}_{h,1}(\Pi_h\sigma - \sigma_h, \eta_r) - ((\text{div} - \beta)\eta_r, P_h u - u_h) = O(h^k),$$

with $k = 1$ or 2 depending on the coefficients in (1) and

$$\bar{a}_{h,3}(\Pi_h\sigma - \sigma_h, \eta_r) - ((\text{div} - \beta)\eta_r, P_h u - u_h) = O(h_{\max}^{k+2}) + \bar{a}_{h,1}(\sigma_h, \eta_r)$$
$$- \bar{a}_{h,3}(\sigma_h, \eta_r).$$

So, $(a_{h,1} - a_{h,3})(\sigma_h, \eta_r)$ is an estimate for the local discretization error. Moreover, this implies that there is a constant $C$ such that

$$\|\Pi_h\sigma - \sigma_h\|_{L^x(\Omega)} + \|P_h u - u_h\|_{L^x(\Omega)} \geq C|\bar{a}_{h,1}(\sigma_h, \eta_r) - \bar{a}_{h,3}(\sigma_h, \eta_r)| + O(h_{\max}^{k+2}).$$

If we assume that

$$\|\Pi_h\sigma - \sigma_h\|_{L^x(\Omega)} + \|P_h u - u_h\|_{L^x(\Omega)} = O(h_{\max}^k),$$

we see that, for $h_{\max}$ small enough,

$$\|\Pi_h\sigma - \sigma_h\|_{L^x(\Omega)} + \|P_h u - u_h\|_{L^x(\Omega)} \geq \frac{1}{2} C|\bar{a}_{h,1}(\sigma_h, \eta_r) - \bar{a}_{h,3}(\sigma_h, \eta_r)|.$$

This provides a lower bound on the global discretization error. We expect the solution for $\bar{a}_{h,3}$ to be two orders of magnitude more accurate than the solution for $\bar{a}_{h,1}$.

## IX. NUMERICAL RESULTS

We consider problem (1) with

$$u = \tanh(\alpha(\sqrt{2}x_1 - ((\sqrt{2} - 1)/2 + x_2))),$$

$$\alpha = 100, \quad \beta_1 = 100, \quad \beta_2 = 100\sqrt{2}, \quad \Gamma_1 = \partial\Omega, \quad g = u|_{\partial\Omega},$$

$$f = -\frac{\text{div}(\text{grad } u + u\beta)}{\alpha}.$$

The smallest upper bound on $h_{\max}$ now follows from the second condition given in Section VIB. This condition reduces to

$$3200 h_{\max} C_r(1 + h_{\max}) < 1/2,$$

so, neglecting the factor $1 + h_{\max}$, for the theorems to apply we would need $h_{\max} \leq 1/(3200C_r)$.

We find the following results for the two discretizations. The directional components of the error vectors for the fluxes were the equal to the accuracy given.

| The $\log_2$ of the errors for $\bar{a}_q = \bar{a}_{h,1}$ | | | |
|---|---|---|---|
| Meshwidth | $\log_2\|P_h u - u_h\|$ | $\log_2\|(\Pi_h\sigma - \sigma_h) \cdot \mathbf{e}_1\|$ | $\log_2\|(\Pi_h\sigma - \sigma_h) \cdot \mathbf{e}_2\|$ |
| 1/4 | $-1.3$ | $-1.3$ | $-1.1$ |
| 1/8 | $-1.8$ | $-1.7$ | $-1.4$ |
| 1/16 | $-2.5$ | $-2.4$ | $-2.1$ |
| 1/32 | $-3.6$ | $-3.4$ | $-3.1$ |
| 1/64 | $-5.1$ | $-5.0$ | $-4.6$ |

| The $\log_2$ of the errors for $\bar{a}_q = \bar{a}_{h,3}$ | | | |
|---|---|---|---|
| Meshwidth | $\log_2\|P_h u - u_h\|$ | $\log_2\|(\Pi_h\sigma - \sigma_h) \cdot \mathbf{e}_1\|$ | $\log_2\|(\Pi_h\sigma - \sigma_h) \cdot \mathbf{e}_2\|$ |
| 1/4 | −2.2 | −2.9 | −2.6 |
| 1/8 | −3.5 | −4.0 | −3.6 |
| 1/16 | −5.4 | −6.1 | −5.6 |
| 1/32 | −8.0 | −9.2 | −8.7 |
| 1/64 | −10.5 | −13.0 | −12.4 |

We see that the order of convergence is indeed higher for the second method. We also see that the difference in order for the fluxes approaches 2. Deviations from the expected order may be caused by the steepness of the solution relative to the mesh.

To test the stability of the low order method, we applied it to problem (1) with

$$u = \tanh(\alpha(\sqrt{2}x_1 - ((\sqrt{2} - 1)/2 + x_2))),$$

with

$$\beta_1 = \alpha, \qquad \beta_2 = \alpha\sqrt{2},$$

for $\alpha = 1000$ and $\alpha = 1000000$, with

$$\Gamma_1 = \partial\Omega, \qquad g = u|_{\partial\Omega}, \qquad f = -\frac{\text{div}(\text{grad } u + u\beta)}{\alpha}.$$

This gave the following results for $\alpha = 1000$.

| The $\log_2$ of the errors for $\bar{a}_q = \bar{a}_{h,1}$ | | | |
|---|---|---|---|
| Meshwidth | $\log_2\|P_h u - u_h\|$ | $\log_2\|(\Pi_h\sigma - \sigma_h) \cdot \mathbf{e}_1\|$ | $\log_2\|(\Pi_h\sigma - \sigma_h) \cdot \mathbf{e}_2\|$ |
| 1/4 | −0.8 | −1.0 | −0.8 |
| 1/8 | −1.0 | −1.2 | −0.8 |
| 1/16 | −1.3 | −1.4 | −1.0 |
| 1/32 | −1.6 | −1.7 | −1.2 |
| 1/64 | −2.0 | −2.0 | −1.5 |

And for $\alpha = 1000000$.

| The $\log_2$ of the errors for $\bar{a}_q = \bar{a}_{h,1}$ | | | |
|---|---|---|---|
| Meshwidth | $\log_2\|P_h u - u_h\|$ | $\log_2\|(\Pi_h\sigma - \sigma_h) \cdot \mathbf{e}_1\|$ | $\log_2\|(\Pi_h\sigma - \sigma_h) \cdot \mathbf{e}_2\|$ |
| 1/4 | −0.8 | −1.0 | −0.8 |
| 1/8 | −1.0 | −1.2 | −0.8 |
| 1/16 | −1.3 | −1.4 | −1.0 |
| 1/32 | −1.5 | −1.6 | −1.2 |
| 1/64 | −1.8 | −1.8 | −1.4 |

The bound given in Section VIB reduces to

$$32\alpha h_{\max} C_r(1 + h_{\max}) < 1/2,$$

so, neglecting the factor $1 + h_{\max}$, for the theorem to apply we need $h_{\max} \leq 1/(32\alpha C_r)$. Evidently, this criterion is not met; nevertheless, convergence occurs, because the scheme reduces to an upwind discretization in cases of large $\alpha$, i.e., for small diffusion constants.

## X. CONCLUSIONS

The Petrov–Galerkin mixed finite element method with exponentially fitted test functions for the fluxes has several nice properties. For instance, just as for a finite volume method, if the true solution $\sigma$ is divergence-free, then the same holds for $\sigma_h$. Furthermore, we have a formal *a priori* error estimate, and, after elimination of $\sigma_h$ by static condensation, the two-dimensional discretization results in an M-matrix for $u_h$. We can extend the method to three dimensions without additional difficulties. Section IX suggests that the scheme with the three-point quadrature rule $\bar{a}_{h,3}$ can serve as a source for *a posteriori* error estimates. To judge the effectiveness of the method for singularly perturbed problems is very difficult. However, the fact that it incorporates exponential fitting, copes well with the exponential solution of the constant coefficient case, and approaches a two-dimensional upwind scheme if the convection goes to infinity suggests that the method based on $\bar{a}_{h,1}$ can be applied to such problems.

## References

1. D. L. Scharfetter and H. K. Gummel, "Large-signal analysis of a silicon read diode oscillator," *IEEE Transactions on Electron Devices*, Vol. ED-16, **1**, 64 (1969).
2. Siegfried Selberherr, *Analysis and Simulation of Semiconductor Devices*, Springer–Verlag, Wien, New York, 1984.
3. Peter A. Markowich, *The Stationary Semiconductor Device Equations*, Springer–Verlag, Wien, New York, 1986.
4. Wolfgang Fichtner, Donald J. Rose, and Randolph E. Bank, "Numerical methods for semiconductor device simulation," *SIAM J. Sci. Stat. Comp.* **4**, 416 (1983).
5. F. Brezzi, L. D. Marini, and P. Pietra, "Mixed exponential fitting schemes for current continuity equations," in *Proceedings of the Sixth International NASECODE Conference*, J. H. Miller, Ed., Press Ltd., Dublin, 1989.
6. P. W. Hemker, *A Numerical Study of Stiff Two-Point Boundary Problems*, Mathematical Centre Tracts, 80, Mathematical Centre, Amsterdam, 1977.
7. J. Douglas, Jr. and J. E. Roberts, "Global estimates for mixed methods for second order elliptic equations," *Math. Comp.* **44**, 39 (1985).
8. Robert A. Adams, *Sobolev Spaces*, Academic Press, New York, 1975.
9. V. Girault and P. Raviart, *Finite Element Methods for Navier–Stokes Equations*, Springer–Verlag, Berlin, 1986.
10. S. J. Polak, C. den Heijer, H. A. Schilders, and P. Markowich, "Semiconductor device modelling for the numerical point of view," *Int. J. Numer. Meth. Eng.* **24**, 763 (1987).
11. Jean E. Roberts and Jean–Marie Thomas, "Mixed and Hybrid Finite Element Methods," RR 737, *INRIA*, Rocquencourt, October 1987.
12. P. A. Raviart and J. M. Thomas, "A mixed finite element for 2-nd order elliptic problems," in *Mathematical Aspects of the Finite Element Method*, Lecture Notes in Mathematics, Springer, **606**, 292 (1977).
13. J. C. Nedelec, "Mixed Finite Elements in $\mathbf{R}^3$," *Numer. Math.* **35**, 315 (1980).
14. P. G. Ciarlet and P. A. Raviart, "General Lagrange and Hermite interpolation in $\mathbf{R}^n$ with applications to finite element methods," *Arch. Rational Mech. Anal.* **46**, 177 (1972).
15. M. Fiedler, *Special Matrices and Their Applications in Numerical Mathematics*, Martinus Nijhoff, Dordrecht, 1986.