

**1991**

R.R.P. van Nooyen

An improved accuracy version of the mixed finite element method  
for a second order elliptic equation

Department of Numerical Mathematics Report NM-R9111 May

CWI is the research institute of the Stichting Mathematisch Centrum, which was founded on February 11, 1946, as a non-profit institution aiming at the promotion of mathematics, computer science, and their applications. It is sponsored by the Dutch Government through the Netherlands organization for scientific research (NWO).

# An Improved Accuracy Version of the Mixed Finite Element Method for a Second Order Elliptic Equation.

R.R.P. van Nooyen

CWI  
P.O. Box 4079, 1009 AB Amsterdam,  
The Netherlands

We discuss a new variant of the mixed finite element method for a second order elliptic problem. By using an appropriate quadrature rule to compute the coefficient matrix, we obtain an improvement in the order of approximation of local averages. We show how the new method can be used to obtain an a posteriori error estimate for a lower order method.

*1980 AMS Subject Classification:* 65N30.

*1990 CR category:* G.1.8 Finite Element Methods.

*Key Words and Phrases:* Mixed Finite Elements, Second Order Elliptic Equations, Quadrature Rules, A Posteriori Error Estimates.

*Note:* This report will be submitted for publication elsewhere.

## 1 Introduction.

In this paper, we describe a modification of the mixed finite element method for a second order elliptic equation. The modified method is based on standard mixed finite elements with lowest order Raviart-Thomas elements on rectangles[1]. To give a short description of our method, we recall, that a mixed finite element formulation of

$$\begin{aligned} - \operatorname{div} a \operatorname{grad} u + cu &= f \text{ on } \Omega, \\ u|_{\partial\Omega} &= g, \end{aligned}$$

can be written as

$$\begin{aligned} (\sigma_h, \tau_h / a) - (\operatorname{div} \tau_h, u_h) &= - \langle g, \tau_h \cdot \mathbf{n}_{\partial\Omega} \rangle \quad \forall \tau_h \in V_h, \\ (\operatorname{div} \sigma_h, t_h) + (cu_h, t_h) &= (f, t_h) \quad \forall t_h \in W_h, \end{aligned}$$

where  $u_h$  is a discrete approximation of  $u$  and  $\sigma_h$  is a discrete approximation of  $-a \operatorname{grad} u$ . In this paper, we show that, if we use a special quadrature rule for the inner product  $(\sigma_h, \tau_h / a)$  and if the coefficient  $a$  is piecewise constant, then the difference between a suitable projection of the continuous solution and the discrete approximation is of order  $\mathcal{O}(h^3)$ . We give numerical evidence that confirms the theoretical result for smooth  $\sigma$ . In section 2, we formulate the boundary value problem to which we apply the modified mixed finite element method. Section 3 describes our mixed finite element discretisation and the quadrature rule for the inner product. There we also give a motivation for the use of the special quadrature rule. We give two other choices for the quadrature rule in section 4. One choice results in the usual scheme for lowest order Raviart-Thomas elements, the other choice corresponds to the use of the trapezoidal rule. We derive an error estimate for the modified version in section 5. In section 6, we use a one dimensional example to illustrate the importance of the ratio  $ch^2/a$  for the usual scheme and our modified scheme. For these methods, the value of this ratio determines whether or not  $u_h$  satisfies a local maximum principle (cf. Polak, Schilders and Couperus[2]) For the scheme based on the trapezoidal rule,  $u_h$  satisfies a local maximum principle for all  $c \geq 0$ . In Section 7, we show numerical results. Section 8 gives an a

posteriori error estimator for the method based on the trapezoidal rule. In the last section, we summarise our results.

## 2 The equation.

We consider a second order elliptic equation with Dirichlet boundary conditions, as given in equation (1),

$$- \operatorname{div} a \operatorname{grad} u + cu = f \text{ on } \Omega, \quad (1a)$$

$$u = g \text{ on } \partial\Omega. \quad (1b)$$

on a rectangle  $\Omega = ]0, L_1[ \times ]0, L_2[$ . We introduce a special notation for  $-a \operatorname{grad} u$ ,

$$\sigma := -a \operatorname{grad} u. \quad (1c)$$

We assume, that there is a finite set of rectangles, the union of which covers  $\Omega$ , such that  $a, c$  are constant on each separate rectangle. We assume that  $a > 0$  and  $c \geq 0$ . We also assume, that  $a, c, f$  and  $g$  are such, that (1) has a unique solution  $u \in C(\Omega)$ , with a  $\sigma$  that is sufficiently smooth for our purposes.

## 3 The discretisation.

In this section, we give a description of our discretisation. We divide  $\Omega$  into rectangular subdomains  $\Omega_{i+\frac{1}{2}, j+\frac{1}{2}}$ , we introduce some notation and we define our test function spaces  $V_h$  and  $W_h$ . We then introduce two projections  $P_h$  and  $\Pi_h$ . Such projections were suggested by Fortin[3] and are used by Raviart and Thomas[1] and Douglas and Roberts[4]. Next, we give the discretisation and discuss the special quadrature rule.

### 3.1. The partitioning of the domain.

We restrict ourselves to subdivisions of the rectangle  $\Omega$ , that can be generated by the Cartesian product of subdivisions of its sides. Let

$$D_1 = \{ 0 = x_{1,0} < x_{1,1} < \dots < x_{1,N_1} = L_1 \}$$

and

$$D_2 = \{ 0 = x_{2,0} < x_{2,1} < \dots < x_{2,N_2} = L_2 \}$$

be partitions such, that  $a$  and  $c$  are constant on the interior of each separate rectangle of the subdivision  $D_1 \times D_2$  of  $\Omega$ . We set

$$h_{1,i+\frac{1}{2}} = x_{1,i+1} - x_{1,i}, \quad (2a)$$

$$h_{2,j+\frac{1}{2}} = x_{2,j+1} - x_{2,j}, \quad (2b)$$

and

$$\Omega_{i+\frac{1}{2}, j+\frac{1}{2}} = \{ (x_1, x_2) \mid x_{1,i} < x_1 < x_{1,i+1}, x_{2,j} < x_2 < x_{2,j+1} \}, \quad (3)$$

$$\Gamma_{i,j+\frac{1}{2}} = \{ (x_1, x_2) \mid x_1 = x_{1,i}, x_{2,j} < x_2 < x_{2,j+1} \}, \quad (4a)$$

$$\Gamma_{i+\frac{1}{2}, j} = \{ (x_1, x_2) \mid x_{1,i} < x_1 < x_{1,i+1}, x_2 = x_{2,j} \}. \quad (4b)$$

### 3.2. The approximation spaces.

We define our approximation spaces for  $\sigma$  and  $u$ , by giving a basis for each space. We then introduce two projections onto the discrete spaces.

For each cell,  $\Omega_{i+\frac{1}{2}, j+\frac{1}{2}}$ , we use the characteristic function  $\chi_{i+\frac{1}{2}, j+\frac{1}{2}}$ ,

$$\chi_{i+\frac{1}{2}, j+\frac{1}{2}} = \delta_{ik} \delta_{jl} \text{ on } \Omega_{k+\frac{1}{2}, l+\frac{1}{2}}, \quad (5)$$

as an element in the set of basis functions for  $W_h$ . For  $V_h$ , we introduce the basis functions  $\eta_{i,j+\frac{1}{2}}$  and  $\eta_{i+\frac{1}{2}, j}$ , where  $\eta_{i,j+\frac{1}{2}}$  is linear in  $x_1$  and constant in  $x_2$  on each cell with

$$\boldsymbol{\eta}_{i,j+\frac{1}{2}} = \delta_{ik} \delta_{jl} \mathbf{e}_1 \text{ on } \Gamma_{k,l+\frac{1}{2}}, \quad (6)$$

for  $i,k=0,1,\dots,N_1, j,l=0,1,\dots,N_2-1$  and  $\boldsymbol{\eta}_{i+\frac{1}{2},j}$  is linear in  $x_2$  and constant in  $x_1$  on each cell with

$$\boldsymbol{\eta}_{i+\frac{1}{2},j} = \delta_{ik} \delta_{jl} \mathbf{e}_2 \text{ on } \Gamma_{k+\frac{1}{2},l}, \quad (7)$$

for  $i,k=0,1,\dots,N_1-1, j,l=0,1,\dots,N_2$ . Here  $\mathbf{e}_1$  and  $\mathbf{e}_2$  are unit vectors in the  $x_1$ - and  $x_2$ -direction respectively.

With these basis functions, we construct  $V_h$  and  $W_h$ ,

$$V_h = \text{Span}(\{ \boldsymbol{\eta}_{i,j+\frac{1}{2}} \mid i=0,1,\dots,N_1, j=0,1,\dots,N_2-1 \} \cup \{ \boldsymbol{\eta}_{i+\frac{1}{2},j} \mid i=0,1,\dots,N_1-1, j=0,1,\dots,N_2 \}), \quad (8)$$

$$W_h = \text{Span}(\{ \chi_{i+\frac{1}{2},j+\frac{1}{2}} \mid i=0,1,\dots,N_1-1, j=0,1,\dots,N_2-1 \}), \quad (9)$$

The product space  $V_h \times W_h$  is the space of lowest order Raviart-Thomas elements. To prepare for the definition of the two projections onto the discrete spaces, we introduce averages over cells and cell boundaries for  $f \in C(\Omega)$ ,

$$P[\Omega_{i+\frac{1}{2},j+\frac{1}{2}}](f) = \frac{1}{\mu(\Omega_{i+\frac{1}{2},j+\frac{1}{2}})} \int_{\Omega_{i+\frac{1}{2},j+\frac{1}{2}}} f \, d\mu, \quad (10)$$

$$P[\Gamma_{i,j+\frac{1}{2}}](f) = \frac{1}{\lambda(\Gamma_{i,j+\frac{1}{2}})} \int_{\Gamma_{i,j+\frac{1}{2}}} f \, d\lambda, \quad (11)$$

$$P[\Gamma_{i+\frac{1}{2},j}](f) = \frac{1}{\lambda(\Gamma_{i+\frac{1}{2},j})} \int_{\Gamma_{i+\frac{1}{2},j}} f \, d\lambda. \quad (12)$$

In the above definitions,  $\lambda$  is the Lebesgue measure on  $\mathbb{R}$  and  $\mu$  is the Lebesgue measure on  $\mathbb{R}^2$ . We define, for all  $u \in L^2(\Omega)$ ,

$$P_h u = P[\Omega_{i+\frac{1}{2},j+\frac{1}{2}}](u) \text{ on } \Omega_{i+\frac{1}{2},j+\frac{1}{2}} \quad \forall i,j, \quad (13)$$

and, for all  $\sigma \in H^1(\Omega)^2$ ,

$$(\Pi_h \sigma)_1 = P[\Gamma_{i,j+\frac{1}{2}}](\sigma_1) \text{ on } \Gamma_{i,j+\frac{1}{2}}, \quad (14a)$$

$$(\Pi_h \sigma)_2 = P[\Gamma_{i+\frac{1}{2},j}](\sigma_2) \text{ on } \Gamma_{i+\frac{1}{2},j} \quad \forall i,j. \quad (14b)$$

The spaces  $V_h$  and  $W_h$  and the projection  $\Pi_h$  were introduced by Raviart and Thomas[1, 3]. The projections have the following special properties,

*Lemma 1.*

$$\forall u \in L^2(\Omega), t_h \in W_h : (u, t_h) = (P_h u, t_h), \quad (15a)$$

$$\forall \sigma \in H^1(\Omega)^2, t_h \in W_h : (\text{div } \sigma, t_h) = (\text{div } \Pi_h \sigma, t_h). \quad (15b)$$

*Proof.*

Equation (15a) follows immediately from the definition of  $P_h$ . Green's formula,

$$\int_{\Omega_{i+\frac{1}{2},j+\frac{1}{2}}} \text{div } \sigma \, d\mu = \int_{\partial\Omega_{i+\frac{1}{2},j+\frac{1}{2}}} \sigma \cdot \mathbf{n}_{\partial\Omega_{i+\frac{1}{2},j+\frac{1}{2}}} \, d\lambda,$$

proves equation (15b)  $\square$

### 3.3. The discretisation scheme.

We first give the discretisation without specifying the quadrature rule. The choice of a quadrature rule is discussed in section 3.4.

We introduce the space

$$V = H(\text{div}, \Omega) := \{ \tau \in L^2(\Omega)^2 \mid \text{div } \tau \in L^2(\Omega) \}, \quad (16)$$

with inner product,

$$(\sigma, \tau)_V = (\sigma, \tau)_{L^1(\Omega)^2} + (\operatorname{div} \sigma, \operatorname{div} \tau)_{L^1(\Omega)}. \quad (17)$$

This space is discussed by Roberts and Thomas[5]. We also introduce

$$W = L^2(\Omega). \quad (18)$$

Note, that  $\Pi_h$  is only defined on  $H^1(\Omega)^2 \subset H(\operatorname{div}, \Omega)$ . In this paper, when we apply  $\Pi_h$  to the  $\sigma$  defined in (1c), the assumption that this  $\sigma$  lies in  $H^1(\Omega)^2$  is included in the condition " $\sigma$  is smooth enough."

We can now write problem (1) in the form:

$$(\sigma, u) \in V \times W,$$

$$\alpha(\sigma, \tau) - (\operatorname{div} \tau, u) = - \langle g, \tau \cdot \mathbf{n}_{\partial\Omega} \rangle \quad \forall \tau \in V, \quad (19a)$$

$$(\operatorname{div} \sigma, t) + (cu, t) = (f, t) \quad \forall t \in W, \quad (19b)$$

where

$$\alpha(\sigma, \tau) := (\sigma, \tau / a) \quad \forall \sigma, \tau \in V. \quad (20)$$

For our discrete problem, we take

$$(\sigma_h, u_h) \in V_h \times W_h,$$

$$\alpha_h(\sigma_h, \tau_h) - (\operatorname{div} \tau_h, u_h) = - \langle g, \tau_h \cdot \mathbf{n}_{\partial\Omega} \rangle \quad \forall \tau_h \in V_h, \quad (21a)$$

$$(\operatorname{div} \sigma_h, t_h) + (cu_h, t_h) = (f, t_h) \quad \forall t_h \in W_h, \quad (21b)$$

where  $\alpha_h$  is a bilinear form on  $V \times V_h$ , that approximates  $\alpha$  and that satisfies

$$\alpha_h(\sigma, \tau_h) = \alpha_h(\Pi_h \sigma, \tau_h). \quad (22)$$

### 3.4. The definition of $\alpha_h$ .

The bilinear form  $\alpha_h$  describes the quadrature rule used to evaluate  $\alpha$ . The idea behind the introduction of a special quadrature rule is the following, if we combine (19), (21) and (22) with the results of lemma 1, we find,

$$\alpha_h(\Pi_h \sigma - \sigma_h, \tau_h) - (\operatorname{div} \tau_h, P_h u - u_h) = \alpha_h(\Pi_h \sigma, \tau_h) - \alpha(\sigma, \tau_h), \quad (23a)$$

$$(\operatorname{div} (\Pi_h \sigma - \sigma_h), t_h) + (ct_h, P_h u - u_h) = 0. \quad (23b)$$

We see, that the only term on the right hand side of this equation is,

$$\alpha_h(\Pi_h \sigma, \tau_h) - \alpha(\sigma, \tau_h). \quad (24)$$

If the discrete problem is uniquely solvable, then it is invertible. In that case this term is a measure for the difference between  $(\Pi_h \sigma, P_h u)$  and  $(\sigma_h, u_h)$ . We now seek to minimise (24). To do this, we construct a special quadrature rule for the evaluation of  $\alpha(\sigma, \tau_h)$  by defining this rule for  $\alpha(\sigma, \eta_1)$  and  $\alpha(\sigma, \eta_2)$ , for each  $\eta_1, \eta_2$  given by (6) and (7). We first introduce the obvious notations,

$$a_{i+\frac{1}{2}, j+\frac{1}{2}} = P[\Omega_{i+\frac{1}{2}, j+\frac{1}{2}}](a),$$

$$\sigma_{1, i, j+\frac{1}{2}} = P[\Gamma_{i, j+\frac{1}{2}}](\sigma_1),$$

$$\sigma_{2, i+\frac{1}{2}, j} = P[\Gamma_{i+\frac{1}{2}, j}](\sigma_2).$$

Our two-dimensional integration rule corresponds to the use of a one-dimensional three-point rule in one direction and exact integration in the other. To simplify the definition of the quadrature rule, we introduce the following functions,

$$A(h, \tilde{h}, L, R) = \frac{hL}{12} + \frac{h\tilde{h}L}{12(h+\tilde{h})} - \frac{\tilde{h}^3 R}{12h(h+\tilde{h})}, \quad (25a)$$

$$B(h, \bar{h}, L, R) = \frac{\bar{h}(\bar{h} + 4h)R}{12h} + \frac{h(\bar{h} + 4h)L}{12\bar{h}}, \quad (25b)$$

$$C(h, \bar{h}, L, R) = A(\bar{h}, h, R, L), \quad (25c)$$

$$D(h, \bar{h}, L, R) = \frac{3hL}{12} + \frac{h\bar{h}L}{12(h+\bar{h})} - \frac{\bar{h}^3 R}{12h(h+\bar{h})}, \quad (25d)$$

$$E(h, \bar{h}, L, R) = \frac{2\bar{h}R}{12} + \frac{\bar{h}^2 R}{12h} + \frac{2hL}{12} + \frac{h^2 L}{12\bar{h}}, \quad (25e)$$

$$F(h, \bar{h}, L, R) = D(\bar{h}, h, R, L). \quad (25f)$$

Where (25a-c) are used in rules for basis functions with their maximum in the interior of  $\Omega$  and (25d-f) are used in rules for basis functions with a their maximum on the boundary of  $\Omega$ .

Now, we define  $\alpha_h$  for all basis functions. We start by defining its action for the  $e_1$  component of  $\sigma$ . We have to distinguish between basis functions with their maximum on the left boundary of  $\Omega$ , (26a), in the interior, (26b), or on the right boundary of  $\Omega$  (26c).

$$\alpha_h(\sigma, \eta_{0, j+\frac{1}{2}}) / h_{2, j+\frac{1}{2}} := \quad (26a)$$

$$D(h_{1, \frac{1}{2}}, h_{1, 1+\frac{1}{2}}, 1/a_{\frac{1}{2}, j+\frac{1}{2}}, 0)\sigma_{1, 0, j+\frac{1}{2}} + E(h_{1, \frac{1}{2}}, h_{1, 1+\frac{1}{2}}, 1/a_{\frac{1}{2}, j+\frac{1}{2}}, 0)\sigma_{1, 1, j+\frac{1}{2}} +$$

$$F(h_{1, \frac{1}{2}}, h_{1, 1+\frac{1}{2}}, 1/a_{\frac{1}{2}, j+\frac{1}{2}}, 0)\sigma_{1, 2, j+\frac{1}{2}},$$

$$\alpha_h(\sigma, \eta_{i, j+\frac{1}{2}}) / h_{2, j+\frac{1}{2}} := \quad (26b)$$

$$A(h_{1, i-\frac{1}{2}}, h_{1, i+\frac{1}{2}}, 1/a_{i-\frac{1}{2}, j+\frac{1}{2}}, 1/a_{i+\frac{1}{2}, j+\frac{1}{2}})\sigma_{1, i-1, j+\frac{1}{2}} +$$

$$B(h_{1, i-\frac{1}{2}}, h_{1, i+\frac{1}{2}}, 1/a_{i-\frac{1}{2}, j+\frac{1}{2}}, 1/a_{i+\frac{1}{2}, j+\frac{1}{2}})\sigma_{1, i, j+\frac{1}{2}} +$$

$$C(h_{1, i-\frac{1}{2}}, h_{1, i+\frac{1}{2}}, 1/a_{i-\frac{1}{2}, j+\frac{1}{2}}, 1/a_{i+\frac{1}{2}, j+\frac{1}{2}})\sigma_{1, i+1, j+\frac{1}{2}},$$

$$\alpha_h(\sigma, \eta_{N_1, j+\frac{1}{2}}) / h_{2, j+\frac{1}{2}} := \quad (26c)$$

$$D(h_{1, N_1-1-\frac{1}{2}}, h_{1, N_1-\frac{1}{2}}, 0, 1/a_{N_1-\frac{1}{2}, j+\frac{1}{2}})\sigma_{1, N_1-2, j+\frac{1}{2}} + E(h_{1, N_1-1-\frac{1}{2}}, h_{1, N_1-\frac{1}{2}}, 0, 1/a_{N_1-\frac{1}{2}, j+\frac{1}{2}})\sigma_{1, N_1-1, j+\frac{1}{2}} +$$

$$F(h_{1, N_1-1-\frac{1}{2}}, h_{1, N_1-\frac{1}{2}}, 0, 1/a_{N_1-\frac{1}{2}, j+\frac{1}{2}})\sigma_{1, N_1, j+\frac{1}{2}},$$

$$\text{for } i=1, 2, \dots, N_1-1, \quad j=0, 1, \dots, N_2-1.$$

Next, we define the rule for basis elements for the  $e_2$  component. Again, we have to distinguish between basis functions with their maximum on the boundary of  $\Omega$ , (26d, 26f), and basis functions with their maximum in the interior (26e).

$$\alpha_h(\sigma, \eta_{i+\frac{1}{2}, 0}) / h_{1, i+\frac{1}{2}} := \quad (26d)$$

$$D(h_{2, \frac{1}{2}}, h_{2, 1+\frac{1}{2}}, 1/a_{i+\frac{1}{2}, \frac{1}{2}}, 0)\sigma_{2, i+\frac{1}{2}, 0} + E(h_{2, \frac{1}{2}}, h_{2, 1+\frac{1}{2}}, 1/a_{i+\frac{1}{2}, \frac{1}{2}}, 0)\sigma_{2, i+\frac{1}{2}, 1} +$$

$$F(h_{2, \frac{1}{2}}, h_{2, 1+\frac{1}{2}}, 1/a_{i+\frac{1}{2}, \frac{1}{2}}, 0)\sigma_{2, i+\frac{1}{2}, 2},$$

$$\alpha_h(\sigma, \eta_{i+\frac{1}{2}, j}) / h_{1, i+\frac{1}{2}} := \quad (26e)$$

$$A(h_{2, j-\frac{1}{2}}, h_{2, j+\frac{1}{2}}, 1/a_{i+\frac{1}{2}, j-\frac{1}{2}}, 1/a_{i+\frac{1}{2}, j+\frac{1}{2}})\sigma_{2, i+\frac{1}{2}, j-1} +$$

$$B(h_{2, j-\frac{1}{2}}, h_{2, j+\frac{1}{2}}, 1/a_{i+\frac{1}{2}, j-\frac{1}{2}}, 1/a_{i+\frac{1}{2}, j+\frac{1}{2}})\sigma_{2, i+\frac{1}{2}, j} +$$

$$C(h_{2, j-\frac{1}{2}}, h_{2, j+\frac{1}{2}}, 1/a_{i+\frac{1}{2}, j-\frac{1}{2}}, 1/a_{i+\frac{1}{2}, j+\frac{1}{2}})\sigma_{2, i+\frac{1}{2}, j+1},$$

$$\alpha_h(\sigma, \eta_{i+\frac{1}{2}, N_2}) / h_{1, i+\frac{1}{2}} := \quad (26f)$$

$$D(h_{2, N_2-1-\frac{1}{2}}, h_{2, N_2-\frac{1}{2}}, 0, 1/a_{i+\frac{1}{2}, N_2-\frac{1}{2}})\sigma_{2, i+\frac{1}{2}, N_2-2} + E(h_{2, N_2-1-\frac{1}{2}}, h_{2, N_2-\frac{1}{2}}, 0, 1/a_{i+\frac{1}{2}, N_2-\frac{1}{2}})\sigma_{2, i+\frac{1}{2}, N_2-1} +$$

$$F(h_{2, N_2-1-\frac{1}{2}}, h_{2, N_2-\frac{1}{2}}, 0, 1/a_{i+\frac{1}{2}, N_2-\frac{1}{2}})\sigma_{2, i+\frac{1}{2}, N_2},$$

$$\text{for } i=0, 1, \dots, N_1-1, \quad \text{for } j=1, 2, \dots, N_2-1, \dots$$

In section 5.1, we show, that for the above choice of coefficients, (24) is  $O(h^3)$ . The use of a three point integration rule means, that we cannot obtain a higher order than this for (24) unless the mesh is uniform and the coefficients are constant on  $\Omega$ , in which case we gain a factor of  $h$  due to symmetry.

#### 4 Other quadrature rules.

If we take different coefficients in our quadrature rule  $\alpha_h$ , we find other variations on the mixed finite element method for lowest order Raviart-Thomas elements.

##### 4.1. Exact evaluation of the form $\alpha$ on test and trial functions.

If we assume piecewise constant coefficients and we use exact integration for the product of test and trial functions, we obtain,

$$A(h, \tilde{h}, L, R) = \frac{hL}{6}, \quad (27a)$$

$$B(h, \tilde{h}, L, R) = \frac{hL}{3} + \frac{\tilde{h}R}{3}, \quad (27b)$$

$$C(h, \tilde{h}, L, R) = A(\tilde{h}, h, R, L), \quad (27c)$$

$$D(h, \tilde{h}, L, R) = \frac{hL}{3}, \quad (27d)$$

$$E(h, \tilde{h}, L, R) = \frac{hL}{6} + \frac{\tilde{h}R}{6} \quad (27e)$$

$$F(h, \tilde{h}, L, R) = D(\tilde{h}, h, R, L), \quad (27f)$$

this choice results in the usual mixed finite element scheme for this choice of test and the trial function spaces.

##### 4.2. Use of the trapezoidal rule.

The use of the trapezoidal rule corresponds to the choice,

$$A(h, \tilde{h}, L, R) = 0, \quad (28a)$$

$$B(h, \tilde{h}, L, R) = \frac{hL}{2} + \frac{\tilde{h}R}{2}, \quad (28b)$$

$$C(h, \tilde{h}, L, R) = A(\tilde{h}, h, R, L), \quad (28c)$$

$$D(h, \tilde{h}, L, R) = \frac{hL}{2}, \quad (28d)$$

$$E(h, \tilde{h}, L, R) = 0 \quad (28e)$$

$$F(h, \tilde{h}, L, R) = D(\tilde{h}, h, R, L). \quad (28f)$$

For this scheme, elimination of  $\sigma_h$  by static condensation is trivial. For  $c \geq 0$ , the resulting matrix is an M-matrix. This implies, that  $u_h$  satisfies a local maximum principle for  $c \geq 0$ . If  $a \equiv 1$  and  $c \equiv 0$  then the matrix after static condensation corresponds to the classical five point finite difference stencil for the Laplace operator.

#### 5 An error estimate.

We derive estimates for  $\|\Pi_h \sigma - \sigma_h\|_{L^2(\Omega)}$  and  $\|P_h u - u_h\|_{L^2(\Omega)}$  under the conditions,

$$c \geq 0 \text{ on } \Omega, \quad (C1)$$

$$\sigma \text{ is smooth enough,} \quad (C2)$$

and



$$A_0(\tau_h, \tau_h) \leq \alpha_h(\tau_h, \tau_h) \leq A_1(\tau_h, \tau_h), \quad (C3)$$

where  $A_0$  and  $A_1$  are positive real numbers, independent of the mesh. To derive error estimates, we need an estimate of the quadrature error, given in section 5.1, and a special norm on  $V_h$ , given in section 5.2. Section 5.3 contains the proof of the error estimate. In section 5.4 we show that condition (C3) is satisfied for a special case.

### 5.1. Error estimates for integration formulas.

We derive an error estimate for our special two dimensional quadrature rule. This rule is based on the interpretation of the values of  $\prod_h \sigma$  on the edges of cells as averages over those edges. Combined with a piecewise constant  $a$  and essentially one-dimensional weight functions, this allows a simple extension of one dimensional integration rules to two dimensions.

To prove this, we combine a special case of Theorem 2 of Bramble and Hilbert[6] with Fubini's theorem[7, 8] and a Sobolev embedding theorem[9]. In lemma 5 we combine these results to give an error estimate. In lemmas 6 and 7 we show that the coefficients given in section 3.4 satisfy the conditions of lemma 5. In lemmas 2, 3 and 4 we formulate the theorems used.

#### Lemma 2.

Let  $\Omega$  be an interval of length  $\rho < \infty$  and let  $1 \leq p < \infty$ . If  $F$  is a linear functional on the Sobolev space  $W^{k,p}(\Omega)$ , which satisfies

$$\exists C > 0 : |F(u)| \leq C \|u\|_{W^{k,p}(\Omega)} \quad \forall u \in W^{k,p}(\Omega), \quad (i)$$

$$F(v) \equiv 0 \quad \forall v \in \{1, x, \dots, x^{k-1}\}, \quad (ii)$$

then

$$\exists \tilde{C} > 0 : |F(u)| \leq \tilde{C} \rho^k \|d^k u / dx^k\|_{L^p(\Omega)}.$$

#### Proof.

This is a special case of Theorem 2 from the article by Bramble and Hilbert[6].

□

#### Lemma 3.

Let  $\Omega_1, \Omega_2$  be bounded intervals in  $\mathbb{R}$ . For  $x \in \Omega_1$ , let  $f[x]$  be the function on  $\Omega_2$  given by  $f[x](y) = f(x, y) \quad \forall y \in \Omega_2$ . If  $f$  is integrable on  $\Omega_1 \times \Omega_2$ , then  $f[x]$  is integrable on  $\Omega_2$  for almost all  $x \in \Omega_1$ ,  $F(x) := \int_{\Omega_2} f[x] d\lambda$  is integrable on  $\Omega_1$  and

$$\int_{\Omega_1 \times \Omega_2} f d\mu = \int_{\Omega_1} F d\lambda.$$

#### Proof.

This is a special case of the theorem of Fubini. [7, 8]

□

We use the Sobolev embedding theorem to give a relation between the maximum norm and the norms on  $W^{2,1}(\Omega)$  and  $W^{2,2}(\Omega)$  if  $\Omega$  is a bounded interval.

#### Lemma 4.

If  $\Omega$  is a bounded interval in  $\mathbb{R}$ , then there are  $C, \tilde{C} > 0$ , such that

$$\|u\|_{L^\infty(\Omega)} \leq C \|u\|_{H^1(\Omega)} \quad \forall u \in H^1(\Omega),$$

$$\|u\|_{L^\infty(\Omega)} \leq \tilde{C} \|u\|_{W^{2,1}(\Omega)} \quad \forall u \in W^{2,1}(\Omega).$$

#### Proof.

Sobolev embedding theorem, see e.g. Girault and Raviart Theorem 1.3[9].

□

The next lemma gives an error estimate for our special two dimensional quadrature rule. To obtain this estimate, we use that our weight functions (i.e. the basis functions  $\eta$ ) are essentially one dimensional. We also use that the values for  $\sigma_h$  can be interpreted as averages over cell edges and that we can define these averages for  $\sigma$ , if  $\sigma$  is smooth enough.

*Lemma 5.*

Let  $\Omega_1, \Omega_2$  and  $\Omega_3$  be bounded intervals in  $\mathbb{R}$ , with  $\Omega_1 \subset \Omega_2$  and  $\rho = \lambda(\Omega_2)$ , the length of  $\Omega_2$ . Furthermore, let  $x_1, x_2, \dots, x_{2k+1} \in \Omega_2$ , let  $w \in L^\infty(\Omega_1)$ ,  $w_1, w_2, \dots, w_{2k+1} \in \mathbb{R}$ , and let  $n \geq 0$ . Set

$$G(u) := \int_{\Omega_1} wu \, d\lambda - \sum_{j=1}^{2k+1} w_j u(x_j) \quad \forall u \in W^{n+1,1}.$$

If

$$G(u) \equiv 0 \quad \forall u \in \{1, x, \dots, x^n\},$$

$f \in C(\overline{\Omega_2 \times \Omega_3})$ ,  $f|_{\Omega_3} \in W^{n+1,1} \quad \forall y \in \Omega_3$ , where we have  $f|_y(x) := f(x, y) \quad \forall x \in \Omega_2$ , and  $\partial^{n+1} f / \partial x^{n+1} \in L^1(\Omega_2 \times \Omega_3)$  then

$$\int_{\Omega_3} G(f) \, d\mu \leq C \rho^{n+1} \|\partial^{n+1} f / \partial x^{n+1}\|_{L^1(\Omega_2 \times \Omega_3)}.$$

*Proof.*

For  $G$ , lemma 4 implies the existence of a  $\tilde{C} > 0$ , such that,

$$G(u) \leq \tilde{C} \|u\|_{W^{2,1}(\Omega_2)} \quad \forall u \in W^{2,1}(\Omega_2).$$

Fubini implies,

$$\left| \int_{\Omega_1 \times \Omega_3} w(x) f(x, y) \, d\mu - \sum_{j=1}^{2k+1} w_j \int_{\Omega_3} f(x_j, y) \, dy \right| = \left| \int_{\Omega_3} G(f) \, dy \right|.$$

We combine this and find, that there exists a  $C > 0$ , such that

$$\left| \int_{\Omega_3} G(f) \, dy \right| \leq C \rho^{n+1} \|\partial^{n+1} f / \partial x^{n+1}\|_{L^1(\Omega_2 \times \Omega_3)}.$$

This follows immediately from lemma 2.

□

In the above lemma,  $G$  corresponds to the error for a one dimensional integration rule. Next, we relate the condition on  $G$  to the coefficients from (24a-f).

*Lemma 6.*

If  $f \in H^3([-h, \tilde{h}])$ ,  $A, B$  and  $C$  are given by (24a-c) and

$$G(f) := L \int_{-h}^0 f(x) \frac{h+x}{h} \, dx + R \int_0^{\tilde{h}} f(x) \frac{\tilde{h}-x}{\tilde{h}} \, dx - \left[ A(h, \tilde{h}, L, R) f(-h) + B(h, \tilde{h}, L, R) f(0) + C(h, \tilde{h}, L, R) f(\tilde{h}) \right], \quad (29)$$

then

$$G(p) \equiv 0 \quad \forall p \in \{1, x, x^2\}.$$

*Proof.*

This can be proved by direct substitution of the appropriate functions in  $G(p)$ .

□

*Lemma 7.*

If  $f \in C^3([-h, \bar{h}])$ ,  $D, E$  and  $F$  are defined by (24d-f) and

$$G(f) := L \int_{-h}^0 f(x) \frac{-x}{h} dx + R \int_0^{\bar{h}} f(x) \frac{x}{\bar{h}} dx - \left[ D(h, \bar{h}, L, R) f(-h) + E(h, \bar{h}, L, R) f(0) + F(h, \bar{h}, L, R) f(\bar{h}) \right], \quad (30)$$

then

$$G(p) \equiv 0 \quad \forall p \in \{ 1, x, x^2 \}.$$

*Proof.*

This is proved as in the previous lemma.

□

Lemma 5 and 6 show, that, we can find a quadrature rule for  $\alpha_h(\cdot, \cdot)$ , that is  $O(h^3)$ . If  $\bar{h} = h$  and  $L = R$ , then we gain an additional order  $h$ ,

*Lemma 8.*

If  $f \in C^4([-h, h])$ ,

$$G(f) := \int_{-h}^0 f(x) \frac{h+x}{h} dx + \int_0^h f(x) \frac{h-x}{h} dx - \left[ \frac{h}{12} f(-h) + \frac{10h}{12} f(0) + \frac{h}{12} f(h) \right] \quad (31)$$

then

$$G(p) \equiv 0 \quad \forall p \in \{ 1, x, x^2, x^3 \}.$$

*Proof.*

Again, this can be proved by calculating  $G(p)$  for the appropriate functions.

□

## 5.2. A special norm on $V_h$ .

The space  $V_h$  is a finite dimensional vector space. Its natural norm is the Euclidean vector norm. For later use, we introduce  $\|\cdot\|_{V_h}$ , a weighted version of the Euclidean vector norm on  $V_h$  and we prove, that this norm is equivalent with the  $L^2(\Omega)$  norm.

If  $\sigma_h \in V_h$  and

$$\sigma_h = \sum_{i=0}^{N_1} \sum_{j=0}^{N_2-1} s_{1,i,j+\frac{1}{2}} \eta_{i,j+\frac{1}{2}} + \sum_{i=0}^{N_1-1} \sum_{j=0}^{N_2} s_{2,i+\frac{1}{2},j} \eta_{i+\frac{1}{2},j}, \quad (32)$$

then we define  $\|\cdot\|_{V_h}$  as,

$$\|\sigma_h\|_{V_h}^2 = \sum_{i=0}^{N_1-1} \sum_{j=0}^{N_2-1} \frac{1}{2} \mu(\Omega_{i+\frac{1}{2},j+\frac{1}{2}}) (s_{1,i,j+\frac{1}{2}}^2 + s_{1,i+1,j+\frac{1}{2}}^2 + s_{2,i+\frac{1}{2},j}^2 + s_{2,i+\frac{1}{2},j+1}^2). \quad (33)$$

*Lemma 9.*

For the  $\sigma_h$  as given in (32),

$$\frac{\|\sigma_h\|_{V_h}^2}{3} \leq \|\sigma_h\|_{L^2(\Omega)}^2 \leq \|\sigma_h\|_{V_h}^2, \quad (34a)$$

and

$$\mu(\Omega_{i+\frac{1}{2},j+\frac{1}{2}}) \|\sigma_h\|_{L^\infty(\Omega_{i+\frac{1}{2},j+\frac{1}{2}})}^2 \leq 2 \|\sigma_h\|_{V_h}^2. \quad (34b)$$

*Proof.*

For both norms, we have

$$\|\sigma_h\|^2 = \|(\sigma_h \cdot \mathbf{e}_1)\mathbf{e}_1\|^2 + \|(\sigma_h \cdot \mathbf{e}_2)\mathbf{e}_2\|^2,$$

so it suffices to prove the inequalities for a single component of  $\sigma_h$ . Furthermore, we know, that

$$\|(\sigma_h \cdot \mathbf{e}_1)\mathbf{e}_1\|_{L^2(\Omega)}^2 = \sum_{i=0}^{N_1-1} \sum_{j=0}^{N_2-1} \|(\sigma_h \cdot \mathbf{e}_1)\mathbf{e}_1\|_{L^2(\Omega_{i+\frac{1}{2},j+\frac{1}{2}})}^2.$$

We compare terms for corresponding cells,

$$\|(\sigma_h \cdot \mathbf{e}_1)\mathbf{e}_1\|_{L^2(\Omega_{i+\frac{1}{2},j+\frac{1}{2}})}^2 = \int_{\Omega_{i+\frac{1}{2},j+\frac{1}{2}}} \left[ s_{1,i,j+\frac{1}{2}} \left[ 1 - \frac{x_1 - x_{1,i}}{h_{1,i+\frac{1}{2}}} \right] + s_{1,i+1,j+\frac{1}{2}} \frac{x_1 - x_{1,i}}{h_{1,i+\frac{1}{2}}} \right]^2.$$

The contribution of

$$\|(\sigma_h \cdot \mathbf{e}_1)\mathbf{e}_1\|_{V_h}^2$$

for this cell is,

$$\frac{1}{2} \mu(\Omega_{i+\frac{1}{2},j+\frac{1}{2}}) (s_{1,i,j+\frac{1}{2}}^2 + s_{1,i+1,j+\frac{1}{2}}^2).$$

The inequalities in (34a) now follow from,

$$\int_0^1 (a\xi + b[1-\xi])^2 d\xi = \frac{a^2 + b^2}{3} + \frac{2ab}{6},$$

and

$$\frac{a^2 + b^2}{6} \leq \frac{a^2 + b^2}{3} + \frac{2ab}{6} \leq \frac{a^2 + b^2}{2}.$$

Inequality (34b) is trivial.

□

### 5.3. The error estimate for the modified method.

In theorem 1, we give an estimate for  $\|\Pi_h \sigma - \sigma_h\|_{L^2(\Omega)}$  and in theorem 2, we give an estimate for  $\|P_h u - u_h\|_{L^2(\Omega)}$ .

*Theorem 1.*

We define,

$$h_1 = \max_i h_{1,i+\frac{1}{2}},$$

$$h_2 = \max_j h_{2,j+\frac{1}{2}},$$

If we assume, that conditions C1 to C3 hold, then

$$\begin{aligned} & \|\Pi_h \sigma - \sigma_h\|_{L^2(\Omega)}^2 + \|\sqrt{c}(P_h u - u_h)\|_{L^2(\Omega)}^2 \leq \\ & K(h_1 + h_2)^3 \max\left(\left\|\frac{\partial^3 \sigma}{\partial x^3}\right\|_{L^\infty(\Omega)}, \left\|\frac{\partial^3 \sigma}{\partial y^3}\right\|_{L^\infty(\Omega)}\right) \\ & \left[ \|\Pi_h \sigma - \sigma_h\|_{L^2(\Omega)} + (h_1 + h_2) \|(\Pi_h \sigma - \sigma_h) \cdot \mathbf{n}_{\partial\Omega}\|_{L^2(\partial\Omega)} \right], \end{aligned} \quad (35a)$$

and

$$\begin{aligned} & \| \Pi_h \sigma - \sigma_h \|_{L^2(\Omega)}^2 + \| \sqrt{c}(P_h u - u_h) \|_{L^2(\Omega)}^2 \leq \\ & K(h_1 + h_2)^3 \max \left( \| \frac{\partial^3 \sigma}{\partial x^3} \|_{L^\infty(\Omega)}, \| \frac{\partial^3 \sigma}{\partial y^3} \|_{L^\infty(\Omega)} \right) \| \Pi_h \sigma - \sigma_h \|_{L^2(\Omega)}. \end{aligned} \quad (35b)$$

*Proof.*

Condition C3 implies, that

$$A_0(\Pi_h \sigma - \sigma_h, \Pi_h \sigma - \sigma_h) \leq \alpha_h(\Pi_h \sigma - \sigma_h, \Pi_h \sigma - \sigma_h).$$

If we set  $\tau = \tau_h$  and  $t = t_h$  in (19) and combine the resulting formulas with (21), we get,

$$\alpha_h(\Pi_h \sigma - \sigma_h, \tau_h) - (\operatorname{div} \tau_h, u - u_h) = \alpha_h(\Pi_h \sigma, \tau_h) - \alpha(\sigma, \tau_h) \quad \forall \tau_h \in \tilde{V}_h, \quad (36a)$$

$$(\operatorname{div}(\sigma - \sigma_h), t_h) + (c(u - u_h), t_h) = 0 \quad \forall t_h \in W_h. \quad (36b)$$

If we take into account (22) and the properties of  $P_h$  and  $\Pi_h$  from lemma 1, then we find,

$$\alpha_h(\Pi_h \sigma - \sigma_h, \tau_h) - (\operatorname{div} \tau_h, P_h u - u_h) = \alpha_h(\sigma, \tau_h) - \alpha(\sigma, \tau_h) \quad \forall \tau_h \in \tilde{V}_h, \quad (37a)$$

$$(\operatorname{div}(\Pi_h \sigma - \sigma_h), t_h) + (c(P_h u - u_h), t_h) = 0 \quad \forall t_h \in W_h. \quad (37b)$$

If we set  $\tau_h = \Pi_h \sigma - \sigma_h$ ,  $t_h = P_h u - u_h$ , then we find

$$\begin{aligned} & \alpha_h(\Pi_h \sigma - \sigma_h, \Pi_h \sigma - \sigma_h) + (c(P_h u - u_h), P_h u - u_h) \\ & = \alpha(\sigma, \Pi_h \sigma - \sigma_h) - \alpha_h(\sigma, \Pi_h \sigma - \sigma_h). \end{aligned}$$

by adding (37b) to (37a).

We introduce

$$K_E = \{ (i, j - 1/2) \mid i = 1, \dots, N_1, j = 1, \dots, N_2 \}, \quad (38a)$$

$$K_N = \{ (i - 1/2, j) \mid i = 1, \dots, N_1, j = 1, \dots, N_2 \}, \quad (38b)$$

$$K_W = \{ (i, j - 1/2) \mid i = 0, \dots, N_1 - 1, j = 1, \dots, N_2 \}, \quad (38c)$$

$$K_S = \{ (i - 1/2, j) \mid i = 1, \dots, N_1, j = 0, \dots, N_2 - 1 \}. \quad (38d)$$

The measure of the support of  $\eta_k$  is denoted by  $\mu(\operatorname{Supp}(\eta_k))$ . We denote the length of the support in the  $\mathbf{e}_\kappa$  direction by  $\lambda_\kappa(\operatorname{Supp}(\eta_k))$ . If  $A$  and  $B$  are sets, we use,

$$A \Delta B = (B - A) \cup (A - B),$$

(the symmetric set difference).

If we combine lemma 5 with lemma 6, lemma 7 and (C3), we find,

$$\begin{aligned} & \alpha(\sigma, \Pi_h \sigma - \sigma_h) - \alpha_h(\sigma, \Pi_h \sigma - \sigma_h) \leq \\ & \sum_{k \in (K_W \cap K_E) \cup (K_W \Delta K_E)} |P[\Gamma_k](\Pi_h \sigma - \sigma_h) \cdot \mathbf{e}_1| (\alpha(\sigma, \eta_k) - \alpha_h(\sigma, \eta_k)) \\ & + \sum_{m \in (K_N \cap K_S) \cup (K_N \Delta K_S)} |P[\Gamma_m](\Pi_h \sigma - \sigma_h) \cdot \mathbf{e}_2| (\alpha(\sigma, \eta_m) - \alpha_h(\sigma, \eta_m)) \leq \\ & C \sum_{k \in (K_W \cap K_E) \cup (K_W \Delta K_E)} |P[\Gamma_k](\Pi_h \sigma - \sigma_h) \cdot \mathbf{e}_1| \mu(\operatorname{Supp}(\eta_k)) \lambda_1(\operatorname{Supp}(\eta_k))^3 \| \frac{\partial^3 \sigma}{\partial x^3} \|_{L^\infty(\Omega)} + \\ & + C \sum_{m \in (K_N \cap K_S) \cup (K_N \Delta K_S)} |P[\Gamma_m](\Pi_h \sigma - \sigma_h) \cdot \mathbf{e}_2| \mu(\operatorname{Supp}(\eta_m)) \lambda_2(\operatorname{Supp}(\eta_m))^3 \| \frac{\partial^3 \sigma}{\partial y^3} \|_{L^\infty(\Omega)}. \end{aligned}$$

From this formula we can derive (35a) and (35b). We start by deriving (35a),

$$\begin{aligned} & \alpha(\sigma, \Pi_h \sigma - \sigma_h) - \alpha_h(\sigma, \Pi_h \sigma - \sigma_h) \leq \\ & C \left[ 2\mu(\Omega) \sum_{k \in K_W \cap K_E} \mu(\operatorname{Supp}(\eta_k)) P[\Gamma_k](\Pi_h \sigma - \sigma_h) \cdot \mathbf{e}_1 \right]^{1/2} 8h_1^3 \| \frac{\partial^3 \sigma}{\partial x^3} \|_{L^\infty(\Omega)} + \end{aligned}$$

$$\begin{aligned}
& C \left[ 2\lambda_1(\Omega) \sum_{k \in (K_w \Delta K_r)} \lambda_1(\Gamma_k) P[\Gamma_k] ((\Pi_h \sigma - \sigma_h) \cdot \mathbf{e}_1)^2 \right]^{\frac{1}{2}} 8h_1^3 h_2 \left\| \frac{\partial^3 \sigma}{\partial x^3} \right\|_{L^\infty(\Omega)} \\
& + C \left[ 2\mu(\Omega) \sum_{m \in K_n \cap K_s} \mu(\text{Supp}(\eta_k)) P[\Gamma_m] ((\Pi_h \sigma - \sigma_h) \cdot \mathbf{e}_2)^2 \right]^{\frac{1}{2}} 8h_2^3 \left\| \frac{\partial^3 \sigma}{\partial y^3} \right\|_{L^\infty(\Omega)} + \\
& C \left[ 2\lambda_2(\Omega) \sum_{m \in (K_n \Delta K_s)} \lambda_2(\Gamma_m) P[\Gamma_m] ((\Pi_h \sigma - \sigma_h) \cdot \mathbf{e}_2)^2 \right]^{\frac{1}{2}} 8h_1 h_2^3 \left\| \frac{\partial^3 \sigma}{\partial y^3} \right\|_{L^\infty(\Omega)} \leq \\
& \tilde{C} \left\| (\Pi_h \sigma - \sigma_h) \right\|_{L^2(\Omega)} (h_1 + h_2)^3 \left[ \left\| \frac{\partial^3 \sigma}{\partial x^3} \right\|_{L^\infty(\Omega)} + \left\| \frac{\partial^3 \sigma}{\partial y^3} \right\|_{L^\infty(\Omega)} \right] + \\
& \tilde{C} \left\| (\Pi_h \sigma - \sigma_h) \cdot \mathbf{n}_{\partial\Omega} \right\|_{L^2(\partial\Omega)} (h_1 + h_2)^4 \left[ \left\| \frac{\partial^3 \sigma}{\partial x^3} \right\|_{L^\infty(\Omega)} + \left\| \frac{\partial^3 \sigma}{\partial y^3} \right\|_{L^\infty(\Omega)} \right].
\end{aligned}$$

Here, we used the equivalence proved in lemma 9. Next, we derive (35b),

$$\begin{aligned}
& \alpha(\sigma, \Pi_h \sigma - \sigma_h) - \alpha_h(\sigma, \Pi_h \sigma - \sigma_h) \leq \\
& C \left[ 2\mu(\Omega) \sum_{k \in K_w \cap K_r} \mu(\text{Supp}(\eta_k)) P[\Gamma_k] ((\Pi_h \sigma - \sigma_h) \cdot \mathbf{e}_1)^2 \right]^{\frac{1}{2}} 8h_1^3 \left\| \frac{\partial^3 \sigma}{\partial x^3} \right\|_{L^\infty(\Omega)} + \\
& C \left[ 2\lambda_1(\Omega) \sum_{k \in (K_w \Delta K_r)} \mu(\text{Supp}(\eta_k)) P[\Gamma_k] ((\Pi_h \sigma - \sigma_h) \cdot \mathbf{e}_1)^2 \right]^{\frac{1}{2}} 8h_1^3 h_2^{\frac{1}{2}} \left\| \frac{\partial^3 \sigma}{\partial x^3} \right\|_{L^\infty(\Omega)} \\
& + C \left[ 2\mu(\Omega) \sum_{m \in K_n \cap K_s} \mu(\text{Supp}(\eta_m)) P[\Gamma_m] ((\Pi_h \sigma - \sigma_h) \cdot \mathbf{e}_2)^2 \right]^{\frac{1}{2}} 8h_2^3 \left\| \frac{\partial^3 \sigma}{\partial y^3} \right\|_{L^\infty(\Omega)} + \\
& C \left[ 2\lambda_2(\Omega) \sum_{m \in (K_n \Delta K_s)} \mu(\text{Supp}(\eta_m)) P[\Gamma_m] ((\Pi_h \sigma - \sigma_h) \cdot \mathbf{e}_2)^2 \right]^{\frac{1}{2}} 8h_1^{\frac{1}{2}} h_2^3 \left\| \frac{\partial^3 \sigma}{\partial x^3} \right\|_{L^\infty(\Omega)} \leq \\
& \tilde{C} \left[ \left\| (\Pi_h \sigma - \sigma_h) \right\|_{L^2(\Omega)} (h_1 + h_2)^3 \left[ \left\| \frac{\partial^3 \sigma}{\partial x^3} \right\|_{L^\infty(\Omega)} + \left\| \frac{\partial^3 \sigma}{\partial y^3} \right\|_{L^\infty(\Omega)} \right] \right].
\end{aligned}$$

Again, we used the equivalence proved in lemma 9.

□

For cells in areas of constant  $a$  and uniform mesh-size, the proof of lemma 8 implies that their contribution to the global error is of order  $h^4$ . If the areas of constant  $a$  and uniform mesh-size are large enough, we treat the cells adjacent to the boundaries of such areas in the same way as the cells adjacent to the boundaries of  $\Omega$ , this results in an  $O(h^{3/2})$  error. If furthermore,

$$\left\| (\Pi_h \sigma - \sigma_h) \cdot \mathbf{n}_{\partial A} \right\|_{L^2(\partial A \cup \partial\Omega)} \leq \left\| \Pi_h \sigma - \sigma_h \right\|_{L^2(\Omega)},$$

where  $\partial A$  is the union of edges between areas of constant  $a$  and uniform mesh-size, then formula (35a) gives us an  $O(h^4)$  error estimate. These effects are seen in our numerical results.

Next, we express  $\|P_h u - u_h\|_{L^1(\Omega)}$  in terms of  $\|\Pi_h \sigma - \sigma_h\|_{L^2(\Omega)}$ .

*Theorem 2.*

Take  $h_1$  and  $h_2$  as in theorem 1. Under the conditions C1, C2 and C3, we have

$$\|P_h u - u_h\|_{L^1(\Omega)} \leq K \left[ \left\| \Pi_h \sigma - \sigma_h \right\|_{L^2(\Omega)} + h_1^3 \left\| \frac{\partial^3 \sigma_1}{\partial x^3} \right\|_{L^\infty(\Omega)} + 2h_1^4 \left\| \frac{\partial^3 \sigma_1}{\partial x^3} \right\|_{L^\infty(\Omega)} \right]. \quad (39)$$

*Proof.*

To obtain this estimate, we examine  $P_h u - u_h$  for each subdomain separately. We use the following relation, which can be obtained from (19) and (21),

$$\alpha(\sigma, \tau_h) - \alpha_h(\sigma_h, \tau_h) - (\operatorname{div} \tau_h, P_h u - u_h) = 0 \quad \forall \tau_h \in V_h.$$

This implies,

$$(\operatorname{div} \tau_h, P_h u - u_h) = \alpha_h(\sigma, \tau_h) - \alpha(\sigma, \tau_h) + \alpha_h(\sigma_h - \Pi_h \sigma, \tau_h) \quad \forall \tau_h \in V_h. \quad (\text{A})$$

We concentrate for the moment on the sub-domain  $\Omega_{i+\frac{1}{2}, j+\frac{1}{2}}$ . We define a special  $\tau_h$ ,

$$\tau_{h,1} = \begin{cases} 0 & \text{on } \Omega_{k+\frac{1}{2}, l+\frac{1}{2}} \text{ if } l < j \text{ or } l > j, \\ 0 & \text{on } \Omega_{k+\frac{1}{2}, j+\frac{1}{2}} \text{ if } k < i, \\ 1 & \text{on } \Omega_{k+\frac{1}{2}, j+\frac{1}{2}} \text{ if } k > i \\ \frac{x_1 - x_{1,i}}{h_{1,i+\frac{1}{2}}} & \text{on } \Omega_{i,j} \end{cases} \quad (40a)$$

$$\tau_{h,2} = 0 \text{ on } \Omega. \quad (40b)$$

Substituting this for  $\tau_h$ , we find,

$$h_{2,j} \|P_h u - u_h\|_{L^\infty(\Omega_{i+\frac{1}{2}, j+\frac{1}{2}})} \leq C \left[ h_{2,j+\frac{1}{2}} \left\| \frac{\partial^3 \sigma_1}{\partial x^3} \right\|_{L^\infty(\Omega)} \left[ \sum_{k=0}^{N_1-1} h_{1,k+\frac{1}{2}}^4 + h_{1,\frac{1}{2}}^4 + h_{1,N_1-\frac{1}{2}}^4 \right] + \sum_{k=0}^{N_1-1} \mu(\Omega_{k+\frac{1}{2}, j+\frac{1}{2}}) \|\Pi_h \sigma - \sigma_h\|_{L^\infty(\Omega_{k+\frac{1}{2}, j+\frac{1}{2}})} \right].$$

The first two terms in this expression correspond with the quadrature error in (A) in the interior and on the edge respectively, the third term corresponds with the remaining term in (A). So,

$$\|P_h u - u_h\|_{L^\infty(\Omega_{i+\frac{1}{2}, j+\frac{1}{2}})} \leq C \left[ (h_1^3 + 2h_1^4) \left\| \frac{\partial^3 \sigma_1}{\partial x^3} \right\|_{L^\infty(\Omega)} + \frac{1}{h_{2,j}} \sum_{k=0}^{N_1-1} \mu(\Omega_{k+\frac{1}{2}, j+\frac{1}{2}}) \|\Pi_h \sigma - \sigma_h\|_{L^\infty(\Omega_{k+\frac{1}{2}, j+\frac{1}{2}})} \right],$$

where we used that  $P_h u - u_h$  is constant on  $\Omega_{i+\frac{1}{2}, j+\frac{1}{2}}$ , Cauchy-Schwartz and (35b). We multiply both sides of this equation by the square root of the area of the cell,

$$\mu(\Omega_{i+\frac{1}{2}, j+\frac{1}{2}})^{\frac{1}{2}} \|P_h u - u_h\|_{L^\infty(\Omega_{i+\frac{1}{2}, j+\frac{1}{2}})} = \|P_h u - u_h\|_{L^2(\Omega_{i+\frac{1}{2}, j+\frac{1}{2}})} \leq \mu(\Omega_{i+\frac{1}{2}, j+\frac{1}{2}})^{\frac{1}{2}} C \left[ (h_1^3 + 2h_1^4) \left\| \frac{\partial^3 \sigma_1}{\partial x^3} \right\|_{L^\infty(\Omega)} + \frac{1}{h_{2,j}} \sum_{k=0}^{N_1-1} \mu(\Omega_{k+\frac{1}{2}, j+\frac{1}{2}}) \|\Pi_h \sigma - \sigma_h\|_{L^\infty(\Omega_{k+\frac{1}{2}, j+\frac{1}{2}})} \right].$$

If we square the left and right hand sides and then sum over  $i$  and  $j$ , we find,

$$\|P_h u - u_h\|_{L^2(\Omega)}^2 \leq K \left[ h_1^6 \left\| \frac{\partial^3 \sigma_1}{\partial x^3} \right\|_{L^\infty(\Omega)}^2 + 2h_1^8 \left\| \frac{\partial^3 \sigma_1}{\partial x^3} \right\|_{L^\infty(\Omega)}^2 + \|\Pi_h \sigma - \sigma_h\|_{L^2(\Omega)}^2 \right].$$

□

Again, if the conditions following the proof of theorem 1 hold, then we gain an additional order of  $h$ , because in that case  $\|\Pi_h \sigma - \sigma_h\|_{L^1(\Omega)}$  is  $O(h^4)$  and we can replace the term  $h_1^3 \left\| \frac{\partial^3 \sigma_1}{\partial x^3} \right\|_{L^\infty(\Omega)}$ , that represents the quadrature error, by  $h_1^4 \left\| \frac{\partial^4 \sigma_1}{\partial x^4} \right\|_{L^\infty(\Omega)}$ .

If, in the above proofs, we replace the explicit expression for the local quadrature error by a more general form, we see, that the order of the error is equal to the order of the quadrature rule used.

#### 5.4. A proof of condition (C3) on a uniform mesh with constant $a$ .

We show that, on a uniform mesh,  $\alpha_h$  satisfies condition (C3) if  $a$  is constant. Without loss of generality we take  $a \equiv 1$ .

*Lemma 10.*

Assume  $a \equiv 1$ . If the mesh is uniform, then

$$\|\sigma_h\|_{V_k}^2 \leq \frac{48}{5} \alpha_h(\sigma_h, \sigma_h) \leq \frac{96}{5} \|\sigma_h\|_{V_k}^2.$$

*Proof.*

If we write  $\sigma_h$  as a linear combination of basis functions  $\boldsymbol{\eta}$ ,

$$\sigma_h = \sum_{m \in K_W \cup K_E} s_{1,m} \boldsymbol{\eta}_m + \sum_{m \in K_N \cup K_S} s_{2,m} \boldsymbol{\eta}_m,$$

then we find,

$$\alpha_h(\sigma_h, \sigma_h) = \alpha_h\left(\sum_{k \in K_W \cup K_E} s_{1,k} \boldsymbol{\eta}_k, \sum_{m \in K_W \cup K_E} s_{1,m} \boldsymbol{\eta}_m\right) + \alpha_h\left(\sum_{k \in K_N \cup K_S} s_{2,k} \boldsymbol{\eta}_k, \sum_{m \in K_N \cup K_S} s_{2,m} \boldsymbol{\eta}_m\right). \quad (41)$$

where  $K_N$  etc. are defined in (38). For the term in (41) corresponding to the  $e_1$  component, we find:

$$\begin{aligned} \alpha_h\left(\sum_{k \in K_W \cup K_E} s_{1,k} \boldsymbol{\eta}_k, \sum_{m \in K_W \cup K_E} s_{1,m} \boldsymbol{\eta}_m\right) = & \\ & h_1 h_2 \sum_{m \in K_W \cap K_E} s_{1,m} \left[ \frac{1}{12} s_{1,m-(1,0)} + \frac{10}{12} s_{1,m} + \frac{1}{12} s_{1,m+(1,0)} \right] + \\ & h_1 h_2 \sum_{m \in K_W - K_E} s_{1,m} \left[ \frac{7}{24} s_{1,m} + \frac{6}{24} s_{1,m+(1,0)} + \frac{-1}{24} s_{1,m+(2,0)} \right] + \\ & h_1 h_2 \sum_{m \in K_E - K_W} s_{1,m} \left[ \frac{-1}{24} s_{1,m-(2,0)} + \frac{6}{24} s_{1,m-(1,0)} + \frac{7}{24} s_{1,m} \right], \end{aligned}$$

where  $m-(1,0) = (i-1, j-\frac{1}{2})$  if  $m = (i, j-\frac{1}{2})$  etc.

Next, we interpret the coefficients  $s_{1,m}$  with  $m \in K_W \cup K_E$  as a vector  $\mathbf{s}$  in  $\mathbb{R}^{(N_1+1)N_2}$ . We introduce the notation  $\mathbf{f}_{1,m}$  for the unit vector along the coordinate axis corresponding to  $s_{1,m}$ . We define the matrix  $A$  by,

$$\mathbf{f}_{1,k}^T A \mathbf{f}_{1,m} = \alpha_h(\boldsymbol{\eta}_{1,k}, \boldsymbol{\eta}_{1,m}).$$

We can write this as follows,

$$\mathbf{s}^T A \mathbf{s} = \frac{1}{2} \mathbf{s}^T (A + A^T) \mathbf{s}.$$

According to the fundamental theorem on symmetric matrices, this implies that all eigenvalues of  $A + A^T$  are real and that,

$$A + A^T = O^T D O,$$

where  $O$  is an orthogonal matrix and  $D$  is a diagonal matrix with as diagonal elements the eigenvalues of  $A$ . Gershgorin's theorem implies that all eigenvalues are larger than

$$\frac{1}{2} h_1 h_2 \left( \frac{14}{24} - \frac{8}{24} - \frac{1}{24} \right) = \frac{5}{48} h_1 h_2,$$

and smaller than

$$h_1 h_2 \left( \frac{10}{12} + \frac{1}{12} + \frac{1}{12} \right) = h_1 h_2.$$

The same reasoning can be applied to the  $e_2$  component of  $\sigma_h$ . We find,



$$\alpha_h(\sigma_h, \sigma_h) = \frac{1}{2} \mathbf{s}^T (A + A^T) \mathbf{s} = \frac{1}{2} \mathbf{s}^T O^T D O \mathbf{s} \geq \frac{5}{48} \|\mathbf{s}\|_{V_h}^2 .$$

□

*Lemma 11.*

For a constant coefficient  $a$ , the bilinear form  $\alpha_h$  satisfies condition (C3).

*Proof.*

This follows immediately from lemma 9 and lemma 10.

□

## 6 The effect of a non-zero $c$ .

We use a one-dimensional example to illustrate the problems associated with a zero order term mentioned in the introduction (cf. Polak, Schilders and Couperus[2] ) The one dimensional problem is studied, because we can easily obtain the discrete system of equations in  $u$ . We see, that, for the quadrature rule given in section 4.1,  $ch^2/a > 6$  results in the loss of the conditions for the local maximum principle for  $u_h$ . For our new quadrature rule, the corresponding bound for satisfying the local maximum principle is  $ch^2/a < 12$ . As any one-dimensional problem can be trivially extended to an example for two dimensions, the same difficulties will appear in two dimensions.

If we write down our modified discretisation in one dimension on a uniform grid with  $a = \epsilon$ ,  $c = 1$ ,  $f = 0$  and  $g(0) = 0$ ,  $g(1) = U$ , then we find the following system of equations:

$$\frac{7h}{24\epsilon} \sigma_0 + \frac{6h}{24\epsilon} \sigma_1 - \frac{h}{24\epsilon} \sigma_2 + u_{\frac{1}{2}} = 0 , \quad (42a.0)$$

$$\frac{h}{12\epsilon} \sigma_{i-1} + \frac{10h}{12\epsilon} \sigma_i + \frac{h}{12\epsilon} \sigma_{i+1} - u_{i-\frac{1}{2}} + u_{i+\frac{1}{2}} = 0 \text{ for } i = 1, 2, \dots, N-1 , \quad (42a.i)$$

$$- \frac{h}{24\epsilon} \sigma_{N-2} + \frac{6h}{24\epsilon} \sigma_{N-1} + \frac{7h}{24\epsilon} \sigma_N - u_{N-\frac{1}{2}} = -U , \quad (42a.N)$$

$$- \sigma_{i-1} + \sigma_i + hu_{i-\frac{1}{2}} = 0 \text{ for } i = 1, 2, \dots, N . \quad (42b)$$

Elimination of  $\sigma$  yields,

$$3u_{\frac{1}{2}} \left(1 + \frac{4h^2}{24\epsilon}\right) - \left(1 - \frac{4h^2}{24\epsilon}\right) u_{1+\frac{1}{2}} = 0 \quad (43a.1)$$

$$- \left(1 - \frac{h^2}{12\epsilon}\right) u_{i-\frac{1}{2}} + 2 \left(1 + \frac{5h^2}{12\epsilon}\right) u_{i+\frac{1}{2}} - \left(1 - \frac{h^2}{12\epsilon}\right) u_{i+1+\frac{1}{2}} = 0 , \quad (43a.i)$$

$$i = 1, 2, \dots, N-2 ,$$

$$- u_{N-1-\frac{1}{2}} \left(1 - \frac{4h^2}{24\epsilon}\right) + 3u_{N-\frac{1}{2}} \left(1 + \frac{4h^2}{24\epsilon}\right) = 2U . \quad (43a.N)$$

We see, that the matrix is always diagonal dominant, but for  $h^2/\epsilon > 12$  it is not an M-matrix.

If we use exact integration for Raviart Thomas mixed finite elements, then we find,

$$3u_{\frac{1}{2}} \left(1 + \frac{h^2}{6\epsilon}\right) - \left(1 - \frac{h^2}{6\epsilon}\right) u_{1+\frac{1}{2}} = 0 \quad (44a.1)$$

$$- \left(1 - \frac{h^2}{6\epsilon}\right) u_{i-\frac{1}{2}} + 2 \left(1 + \frac{h^2}{3\epsilon}\right) u_{i+\frac{1}{2}} - \left(1 - \frac{h^2}{6\epsilon}\right) u_{i+1+\frac{1}{2}} = 0 , \quad (44a.i)$$

$$i = 1, 2, \dots, N-2 ,$$

$$- u_{N-1-\frac{1}{2}} \left(1 - \frac{h^2}{6\epsilon}\right) + 3u_{N-\frac{1}{2}} \left(1 + \frac{h^2}{6\epsilon}\right) = 2U . \quad (44a.N)$$

Here we see, that there is no qualitative difference in sensitivity to the ratio  $\frac{h^2}{\epsilon}$  between our method and the standard method. However, for the trapezoidal rule we find:

$$3u_{\frac{1}{2}}(1 + \frac{h^2}{3\epsilon}) - u_{1+\frac{1}{2}} = 0 \quad (45a.1)$$

$$-u_{i-\frac{1}{2}} + 2(1 + \frac{h^2}{2\epsilon})u_{i+\frac{1}{2}} - u_{i+1+\frac{1}{2}} = 0 \text{ for } i=1,2,\dots,N-2, \quad (45a.i)$$

$$-u_{N-1-\frac{1}{2}} + 3u_{N-\frac{1}{2}}(1 + \frac{h^2}{3\epsilon}) = 2U. \quad (45a.N)$$

In this case, we do get an M-matrix.

We recall from section 5, that the accuracy of a method is determined by the accuracy of the quadrature rule used in  $\alpha_h$ . If  $\sigma$  is sufficiently smooth, then we find the following orders for the above schemes,  $O(h^{3/2})$  for (42) ( $O(h^4)$  if the error is not concentrated at the edges),  $O(h^2)$  for (45) and for (44). The latter result may seem strange, because this scheme is based on exact integration of products of test and trial functions. However, by inspection of the formulas, we see that the need to integrate products of continuous piecewise linear functions results in coefficients, that are not optimal for approximate integration of products of smooth functions and continuous, piecewise linear functions.

In scheme (42) and (44), we find the same equations for boundary cells. The equations for boundary cells in (45) however, are different. As (42) is  $O(h^{3/2})$  accurate, the equations for the boundary given by this scheme are more accurate than those given by (45). So, on the same mesh, we expect the error in the boundary cells for scheme (45) to be larger than for scheme (44), but we expect to find the same order behaviour for both schemes. Our experiments confirm this expectation.

## 7 Numerical experiments.

This section gives numerical results for problem (1) on a uniform mesh. We take  $c \equiv 0$ ,  $\Omega =$  the unit square and  $f, g$  such that

$$u = \frac{(\exp(x - \frac{1}{2}) - 1)(\exp(y - \frac{1}{2}) - 1)}{a},$$

is the solution of the continuous problem. First we give results for  $a=1$  on the unit square, then we divide the unit square into four smaller squares and give results for a discontinuous coefficient  $a$ ,  $a=1$  in the lower left square,  $a=10$  in the upper left square,  $a=100$  in the lower right square and  $a=1000$  in the upper right square (Figure 1).

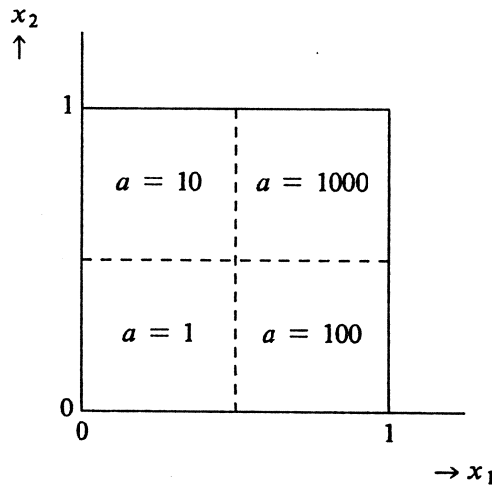


Figure 1.

For  $u_h$ , the size of the error is expressed in the  $L^2(\Omega)$  norm. For  $\sigma_h$ , the size of the error is expressed as the Euclidean norm in the space of vectors of coefficients of the  $\boldsymbol{\eta}$  basis vectors, scaled by the square root of the area of one cell.

We give results for the discretisation from section 3 and the two discretisations from section 4. We indicate the quadrature rule used in the discretisation by roman numbers, I denotes the quadrature given in section 3.4, number II denotes exact quadrature (section 4.1) and discretisation III denotes the trapezoidal rule (section 4.2).

$h$	$\log_2 \ P_h u - u_h\ _E$			$\log_2 \ \Pi_h \sigma - \sigma_h\ _E$		
	I	II	III	I	II	III
1/2	-13.03	-13.03	-7.22	-9.18	-7.43	-5.13
1/4	-16.32	-14.88	-9.13	-12.67	-9.25	-6.42
1/8	-20.06	-16.81	-11.04	-16.41	-11.21	-8.05
1/16	-23.95	-18.79	-12.99	-20.25	-13.21	-9.83
1/32	-27.90	-20.78	-14.97	-24.14	-15.21	-11.67
1/64	-31.87	-22.78	-16.97	-28.07	-17.21	-13.56
1/128	-35.86	-24.78	-18.97	-32.01	-19.21	-15.47
1/256	-39.86	-26.78	-20.97	-35.96	-21.21	-17.39

**Table 1**

*Errors for the three methods for the constant coefficient case.*

$h$	$\log_2 \ P_h u - u_h\ _E$			$\log_2 \ \Pi_h \sigma - \sigma_h\ _E$		
	I	II	III	I	II	III
1/2	-16.23	-16.23	-8.23	-9.26	-7.59	-4.91
1/4	-17.89	-16.94	-10.22	-12.53	-9.31	-6.12
1/8	-21.75	-18.72	-12.20	-16.24	-11.26	-7.71
1/16	-25.69	-20.67	-14.18	-20.06	-13.25	-9.47
1/32	-29.67	-22.65	-16.17	-23.94	-15.25	-11.31
1/64	-33.67	-24.65	-18.17	-27.85	-17.25	-13.18
1/128	-37.67	-26.65	-20.17	-31.78	-19.25	-15.08
1/256	-41.67	-28.65	-22.17	-35.72	-21.25	-17.00

**Table 2**

*Errors for the three methods for the discontinuous coefficient case.*

Starting at  $h = 1/8$ , we see, for case I, convergence of order 4 as predicted in section 5.2 for a uniform mesh and large areas with constant coefficients. The other schemes show second order behaviour. We recall, that the error analysis in section 5 shows that the accuracy of a method is determined by the accuracy of the quadrature rule  $\alpha_h$  applied to  $\sigma$ . Our  $\sigma$  is smooth, so we indeed expect the following orders for the above schemes,  $O(h^4)$  for (I),  $O(h^2)$  for (II),  $O(h^2)$  for (III).

## 8 An a-posteriori error estimate.

We see that there is a difference in order of accuracy between our special method, given in section 3.4 and the method based on the use of the trapezoidal rule, given in section 4.2. This suggests that the special scheme may be used to obtain an a-posteriori estimate of the error in the solution of the trapezoidal scheme.

In this section, we shall use the following notation,  $\alpha_{h,3}$  is the bilinear form we obtain if we use the three point rule given in section 3.4 to evaluate  $\alpha_h$  and  $\alpha_{h,1}$  is the bilinear form we obtain if we use the trapezoidal rule given in section 4.2. Furthermore, let  $(\sigma, u)$  be the solution of problem (19), let  $(\sigma_h, u_h)$  be the solution of the discretisation (21) given in section 3.3 with  $\alpha_h = \alpha_{h,1}$  and let  $(\tilde{\sigma}_h, \tilde{u}_h)$  be the solution of the same discretisation, with  $\alpha_h = \alpha_{h,3}$ .

The simplest way to obtain an a-posteriori error estimate is to solve both schemes. Given the solution of both schemes, we can obtain estimates for

$$\| \Pi_h \sigma - \sigma_h \|_{H(\text{div}, \Omega)},$$

and

$$\| P_h u - u_h \|_{L^2(\Omega)},$$

as follows, we insert an extra term in the above expressions and use the triangle inequality to find,

$$\begin{aligned} \| \tilde{\sigma}_h - \sigma_h \|_{H(\text{div}, \Omega)} - \| \Pi_h \sigma - \tilde{\sigma}_h \|_{H(\text{div}, \Omega)} &\leq \| \Pi_h \sigma - \sigma_h \|_{H(\text{div}, \Omega)} \leq \\ &\| \tilde{\sigma}_h - \sigma_h \|_{H(\text{div}, \Omega)} + \| \Pi_h \sigma - \tilde{\sigma}_h \|_{H(\text{div}, \Omega)}, \end{aligned}$$

and

$$\begin{aligned} \| \tilde{u}_h - u_h \|_{L^2(\Omega)} - \| P_h u - \tilde{u}_h \|_{L^2(\Omega)} &\leq \| P_h u - u_h \|_{L^2(\Omega)} \leq \\ &\| \tilde{u}_h - u_h \|_{L^2(\Omega)} + \| P_h u - \tilde{u}_h \|_{L^2(\Omega)}. \end{aligned}$$

Next, we assume that  $\sigma$  is sufficiently smooth and we recall that

$$\| \Pi_h \sigma - \sigma_h \|_{H(\text{div}, \Omega)} + \| P_h u - u_h \|_{L^2(\Omega)} = \mathcal{O}(h^k).$$

and

$$\| \Pi_h \sigma - \tilde{\sigma}_h \|_{H(\text{div}, \Omega)} + \| P_h u - \tilde{u}_h \|_{L^2(\Omega)} = \mathcal{O}(h^{l+2}),$$

where  $k, l = 2$  if the mesh is uniform and  $a$  is constant and otherwise  $k = 1$  or  $2, l = 1$  or  $2$  depending on the mesh and  $a$ . This implies, that

$$\| \Pi_h \sigma - \sigma_h \|_{H(\text{div}, \Omega)} = (1 + \mathcal{O}(h)) \| \tilde{\sigma}_h - \sigma_h \|_{H(\text{div}, \Omega)},$$

and

$$\| P_h u - u_h \|_{L^2(\Omega)} = (1 + \mathcal{O}(h)) \| \tilde{u}_h - u_h \|_{L^2(\Omega)}.$$

where  $h$  is the maximum cell diameter of the mesh.

## 9 Conclusions.

For equation (1), we have increased the accuracy of the mixed finite element approximation of  $(\Pi_h \sigma, P_h u)$  by introducing a particular quadrature rule for  $\alpha(\sigma, \tau_h)$ . This leads to a scheme, that has the same complexity as standard mixed finite elements for lowest order Raviart-Thomas elements, but that is of  $\mathcal{O}(h^3)$  instead of  $\mathcal{O}(h^2)$  if  $\sigma$  is sufficiently smooth. This behaviour is confirmed by numerical experiments.

In section 8, we show that this difference in order can be used to give an a posteriori error estimator for the less accurate version.

If we compare the usual method (section 4.1) with the other two methods, we see, that the only advantage of the method given in section 4.1 over the method that uses the trapezoidal rule (section 4.2) is a better treatment of boundary cells (see the discussion in section 6). The only advantage of the method given in section 4.1 over our modified method is, that the method from section 4.1 may give exact results for less smooth solutions, viz. for solutions with  $\sigma \in V_h$ .

To decide whether to use the method based on the trapezoidal rule or our modified method, we must weigh the advantage of a simpler matrix, that reduces to an M-matrix for  $u_h$  for all  $c \geq 0$ , against the loss of accuracy. The numerical experiments show the loss of accuracy to be considerable for smooth  $\sigma$ . So, only if it is known, that the combination of  $a, c$  and  $h$  may lead to instability (for instance if  $ch^2/a \geq 1$ ), or if  $\sigma$  is not smooth enough, is it more efficient to use the method based on the trapezoidal rule. In all other cases our modified method would be the better choice.

The choice between our method and the method discussed in section 4.1 is simple. Both methods are equally sensitive to a zero order term. Both methods also have the same sparsity pattern in their matrices, so they roughly need the same amount of work to solve. As the method

section in 4.1 is of lower order than the modified method, the modified method is more efficient if we look at accuracy obtained versus complexity.

## References

1. P. A. Raviart and J. M. Thomas, "A mixed finite element method for 2-nd order elliptic problems," in *Mathematical aspects of the finite element method*, Springer (1977).
2. S. J. Polak, W. H. A. Schilders, and H. D. Couperus, "A finite element method with current conservation," pp. 453-462 in *Simulation of semiconductor devices and processes*, ed. M. Rudan, Tecnoprint, Bologna (1988).
3. M. Fortin, "An analysis of the convergence of mixed finite element methods," *RAIRO Numerical Analysis* **11**(4), pp. 341-354 (1977).
4. J. Douglas, Jr. and J. E. Roberts, "Global estimates for mixed methods for second order elliptic equations," *Mathematics of computation* **44**(169), pp. 39-52 (1985).
5. Jean E. Roberts and Jean-Marie Thomas, "Mixed and Hybrid Finite Element Methods," RR 737, INRIA, Rocquencourt (October 1987).
6. J. H. Bramble and S. R. Hilbert, "Estimation of linear functionals on Sobolev spaces with application to Fourier transforms and spline interpolation," *SIAM J. Numer. Anal.* **7**(1) (1970).
7. P. R. Halmos, *Measure Theory*, Springer Verlag (1974).
8. H. L. Royden, *Real Analysis, second edition*, MacMillan Company (1963).
9. V. Girault and P. Raviart, *Finite Element Methods for Navier-Stokes Equations*, Springer-Verlag (1986).

