

1 **A MONOTONE Q^1 FINITE ELEMENT METHOD FOR**
2 **ANISOTROPIC ELLIPTIC EQUATIONS**

3 HAO LI * AND XIANGXIONG ZHANG †

4 **Abstract.** We construct a monotone continuous Q^1 finite element method on the uniform mesh
5 for the anisotropic diffusion problem with a diagonally dominant diffusion coefficient matrix. The
6 monotonicity implies the discrete maximum principle. Convergence of the new scheme is rigorously
7 proven. On quadrilateral meshes, the matrix coefficient conditions translate into specific a mesh
8 constraint.

9 **Key words.** Inverse positivity, Q^1 finite element method, monotonicity, discrete maximum
10 principle, anisotropic diffusion

11 **AMS subject classifications.** 65N30, 65N15, 65N12

12 **1. Introduction.**

13 **1.1. Monotonicity and discrete maximum principle.** Consider solving the
14 following elliptic equation on $\Omega = (0, 1)^2$ with Dirichlet boundary conditions:

15 (1.1)
$$\begin{aligned} \mathcal{L}u &\equiv -\nabla \cdot (\mathbf{a}\nabla u) + cu = f && \text{on } \Omega, \\ u &= g && \text{on } \partial\Omega, \end{aligned}$$

where the diffusion matrix $\mathbf{a}(\mathbf{x}) \in \mathbb{R}^{2 \times 2}$, $c(\mathbf{x})$, $f(\mathbf{x})$ and $g(\mathbf{x})$ are sufficiently smooth functions over $\bar{\Omega}$ or $\partial\Omega$. We assume that $\forall \mathbf{x} \in \Omega$, $\mathbf{a}(\mathbf{x})$ is symmetric and uniformly positive definite on Ω . In the literature, (1.1) is called a heterogeneous anisotropic diffusion problem when the eigenvalues of $\mathbf{a}(\mathbf{x})$ are unequal and vary over on Ω . For a smooth function $u \in C^2(\Omega) \cap C(\bar{\Omega})$, a maximum principle holds [7]:

$$\mathcal{L}u \leq 0 \quad \text{on } \Omega \quad \implies \quad \max_{\Omega} u \leq \max \left\{ 0, \max_{\partial\Omega} u \right\}.$$

16 In particular,

17 (1.2)
$$\mathcal{L}u = 0 \text{ in } \Omega \implies |u(x_1, x_2)| \leq \max_{\partial\Omega} |u|, \quad \forall (x_1, x_2) \in \Omega.$$

18 For simplicity, we only consider the homogeneous Dirichlet boundary condition,
19 i.e. $g = 0$. The anisotropic diffusion problem (1.1) arises from various areas of
20 science and engineering, including plasma physics, Lagrangian hydrodynamics, and
21 image processing. To avoid spurious oscillations or non-physical numerical solution,
22 it is desired to have numerical schemes to satisfy (1.2) in the discrete sense. We
23 are interested in a linear approximation to \mathcal{L} which can be represented as a matrix
24 L_h . The matrix L_h is called monotone if its inverse only has nonnegative entries,
25 i.e., $L_h^{-1} \geq 0$. Monotonicity of the scheme is a sufficient condition for the discrete
26 maximum principle and has various applications especially for parabolic problems,
27 see [1, 33, 14, 9, 31, 21, 5, 6, 22, 21, 13, 16].

*lihao199408@gmail.com,

†Department of Mathematics, Purdue University, 150 N. University Street, West Lafayette, IN 47907-2067 (zhan1966@purdue.edu).

28 **1.2. Monotone schemes for anisotropic diffusion equations.** Monotone
 29 (or positive-type in some literature) numerical methods for problem (1.1) have received
 30 considerable attention, e.g., see [11, 17, 18, 19, 20, 25, 34, 30, 12, 2, 27]. The major
 31 efforts of studying linear monotone schemes take advantage of M -matrix (see [29] for
 32 the definition), either by showing the coefficient matrix is M -matrix directly or the
 33 coefficient matrix can be factorized into product of M -matrices. In the following, we
 34 call a numerical scheme satisfying M -matrix property if the corresponding coefficient
 35 matrix is a M -matrix.

36 By factorizing the stiffness matrix into a product of M -matrices, the monotonicity
 37 can still be ensured. For a nine-point scheme on a two-dimensional quadrilateral grid,
 38 the matrix condition for monotonicity with specific splitting strategy in [28] aligns
 39 with the Lorenz's condition presented in [23, 14]. The difference is in [23, 14], only
 40 the existence of the factorization was proved while in [28] the authors found the exact
 41 matrix factorization.

42 In [26], it is proved that a monotone finite difference scheme exists for any linear
 43 second-order elliptic problem on fine enough uniform mesh and a finite difference
 44 method with fixed stencil for all the problems satisfying the M -matrix property does
 45 not exist. With nonnegative directional splittings, [32, 8, 27] propose to construct
 46 finite difference schemes for elliptic operators in the nondivergence form and diver-
 47 gence form. Particularly in [27], it is shown that a monotone scheme satisfying the
 48 M -matrix property can be constructed for continuous diffusion matrix for sufficiently
 49 fine mesh and sufficiently large finite difference stencil.

50 In [17], for the P^1 finite elements in two and three dimensions, the author gen-
 51 eralized the well known non-obtuse angle condition for anisotropic diffusion problem
 52 in the sense to have the dihedral angles of all mesh elements, measured in a metric
 53 depending on $\mathbf{a}(\mathbf{x})$, be non-obtuse. It reduces to the non-obtuse angle condition for
 54 isotropic diffusion matrices when $\mathbf{a}(\mathbf{x}) = \alpha(\mathbf{x})\mathbb{I}$. The formulation was also utilized in
 55 [17] for the construction of the so called M -uniform meshes on which the numerical
 56 scheme is monotone. The approach to show monotonicity in [17] is to write the global
 57 matrix as the sum of local contributions. In [10], the Delaunay condition is extended
 58 to anisotropic diffusion problems through a refined analysis studying the whole stiff-
 59 ness matrix for the two-dimensional situation. The analysis of [17] was extended to
 60 the anisotropic diffusion-convection-reaction problems in [24].

61 For the Q^1 finite elements, research on monotonicity has predominantly been
 62 focused on meshes whose cells are rectangular blocks. For the two-dimensional Poisson
 63 equation, it was noted in [3] that the M -matrix property is violated when the aspect
 64 ratio, i.e. the ratio between the length of the longer edge and the shorter edge of
 65 the cell, becomes excessively large. Then the discrete maximum principle is not
 66 guaranteed.

67 **1.3. Contributions and organization of the paper.** It is well known that
 68 the second-order accurate linear schemes, such as mixed finite element and multi-
 69 point flux approximation, do not always satisfy monotonicity for distorted meshes or
 70 with high anisotropy ratio. In this paper, we construct a monotone Q^1 finite element
 71 method for solving the equation (1.1), which is second-order accurate for function
 72 values.

73 To analyze the monotonicity of the stiffness matrix, we approximate integrals
 74 with a specific quadrature rule, particularly, the linear combination of the trapezoid
 75 rule and midpoint rule. We demonstrate that a continuous Q^1 finite element method
 76 with the specific quadrature rule, when applied to the anisotropic diffusion problem

77 on a uniform mesh, ensures monotonicity for the problem with a diagonally domi-
 78 nant diffusion coefficient matrix. The method is linear, second-order accurate. The
 79 convergence of the function values for this method is also proven. The coefficient
 80 constraints become mesh constraints when this Q^1 finite element method is used on
 81 general quadrilateral meshes.

82 The paper is organized as follows. In Section 2, we introduce the notations and
 83 review standard quadrature estimates. In Section 3, we derive the Q^1 scheme for
 84 anisotropic diffusion equation with Dirichlet boundary condition and derive the coef-
 85 ficient constraints for the stiffness matrix to be an M -matrix. In Section 4, we prove
 86 the convergence of function values. In Section 5, we discuss the extension to general
 87 quadrilateral meshes. Numerical results are given in Section 6.

88 2. Preliminaries.

89 **2.1. Notation and tools.** We list the tools and notation as follows.

- 90 • For the problem dimension d , though we only consider the case $d = 2$, some-
 91 times we keep the general notation d to illustrate how the results are influ-
 92 enced by the dimension.
- 93 • For the Q^1 finite element space, i.e., tensor product of linear polynomials, the
 94 local space is defined on a reference cell \hat{K} , e.g., $\hat{K} = [0, 1]^2$. Then, the finite
 95 element space on a physical mesh cell e is given by the reference map from
 \hat{K} to e . The reference element \hat{K} is as Figure 1.

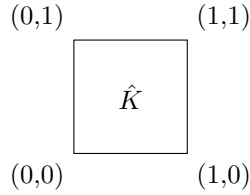


FIG. 1. *The reference element.*

96 On a reference element \hat{K} , we have the Lagrangian basis $\hat{\phi}_{0,0}, \hat{\phi}_{0,1}, \hat{\phi}_{1,1}, \hat{\phi}_{1,0}$
 97 as
 98

$$(2.1) \quad \hat{\phi}_{0,0} = (1-\hat{x}_1)(1-\hat{x}_2), \quad \hat{\phi}_{0,1} = (1-\hat{x}_1)\hat{x}_2, \quad \hat{\phi}_{1,1} = \hat{x}_1\hat{x}_2, \quad \hat{\phi}_{1,0} = \hat{x}_1(1-\hat{x}_2).$$

- 100 • We will use $\hat{\cdot}$ for a function to emphasize the function is defined on or trans-
 101 formed to the reference element \hat{K} from a physical mesh element.
- For a quadrilateral element e , we assume \mathbf{F}_e is the bilinear mapping such that
 $\mathbf{F}_e(\hat{K}) = e$. Let $\mathbf{c}_{i,j}$, $i, j = 0, 1$ be the vertices of the quadrilateral element e .
 The mapping \mathbf{F}_e can be written as

$$\mathbf{F}_e = \sum_{\ell=0}^1 \sum_{m=0}^1 \mathbf{c}_{\ell,m} \hat{\phi}_{\ell,m}.$$

- 102 • $Q^1(\hat{K}) = \left\{ p(\mathbf{x}) = \sum_{i=0}^1 \sum_{j=0}^1 p_{ij} \hat{\phi}_{i,j}(\hat{\mathbf{x}}), \quad \hat{\mathbf{x}} \in \hat{K} \right\}$ is the set of Q^1 poly-
 103 nomials on the reference element \hat{K} .
- 104 • $Q^1(e) = \left\{ v_h \in H^1(e) : v_h \circ \mathbf{F}_j \in Q_1(\hat{K}) \right\}$ is the set of Q^1 polynomials on an
 105 element e .

- 106 • $V^h = \{p(\mathbf{x}) \in H^1(\Omega_h) : p|_e \in Q^1(e), \quad \forall e \in \Omega_h\}$ denotes the continuous Q^1
 107 finite element space on Ω_h .
 108 • $V_0^h = \{v_h \in V^h : v_h = 0 \quad \text{on} \quad \partial\Omega\}$
 109 • Let $(f, v)_e$ denote the inner product in $L^2(e)$ and (f, v) denote the inner
 110 product in $L^2(\Omega)$:

$$111 \quad (f, v)_e = \int_e f v \, d\mathbf{x}, \quad (f, v) = \int_{\Omega} f v \, d\mathbf{x} = \sum_e (f, v)_e.$$

- 112 • Let $\langle f, v \rangle_{e,h}$ denote the approximation to $(f, v)_e$ by the mixed quadrature
 113 defined in (2.7) over element e with some specified quadrature parameter and
 114 $\langle f, v \rangle_h$ denotes the approximation to (f, v) by

$$115 \quad \langle f, v \rangle_h = \sum_e \langle f, v \rangle_{e,h}.$$

- 116 • Let $E(f)$ denote the quadrature error for integrating $f(\hat{\mathbf{x}})$ on element e . Let
 117 $\hat{E}(\hat{f})$ denote the quadrature error for integrating $\hat{f}(\hat{\mathbf{x}}) = f(\mathbf{F}_e(\hat{\mathbf{x}}))$ on the
 118 reference element \hat{K} . Then $E(f) = h^d \hat{E}(\hat{f})$ on uniform rectangular mesh
 119 with mesh size h .
 120 • The norm and semi-norms for $W^{k,p}(\Omega)$ and $1 \leq p < +\infty$, with standard
 121 modification for $p = +\infty$:

$$122 \quad \|u\|_{k,p,\Omega} = \left(\sum_{i+j \leq k} \iint_{\Omega} |\partial_{x_1}^i \partial_{x_2}^j u(x_1, x_2)|^p \, d\mathbf{x} \right)^{1/p},$$

$$123 \quad |u|_{k,p,\Omega} = \left(\sum_{i+j=k} \iint_{\Omega} |\partial_{x_1}^i \partial_{x_2}^j u(x_1, x_2)|^p \, d\mathbf{x} \right)^{1/p},$$

$$124 \quad [u]_{k,p,\Omega} = \left(\iint_{\Omega} |\partial_{x_1}^k u(x_1, x_2)|^p \, d\mathbf{x} + \iint_{\Omega} |\partial_{x_2}^k u(x_1, x_2)|^p \, d\mathbf{x} \right)^{1/p}.$$

- 126 • In the special case where $\omega = \Omega$, we drop the subscript, i.e. $(\cdot, \cdot) := (\cdot, \cdot)_{\Omega}$
 127 and $\|\cdot\| := \|\cdot\|_{\Omega}$.
 128 • For any $v_h \in V^h$, $1 \leq p < +\infty$ and $k \geq 1$, we will abuse the notation to
 129 denote the broken Sobolev norm and semi-norms by the following symbols

$$130 \quad \|v_h\|_{k,p,\Omega} := \left(\sum_e \|v_h\|_{k,p,e}^p \right)^{\frac{1}{p}},$$

$$131 \quad |v_h|_{k,p,\Omega} := \left(\sum_e |v_h|_{k,p,e}^p \right)^{\frac{1}{p}},$$

$$132 \quad [v_h]_{k,p,\Omega} := \left(\sum_e [v_h]_{k,p,e}^p \right)^{\frac{1}{p}}.$$

- 134 • For simplicity, sometimes we may use $\|u\|_{k,\Omega}$, $|u|_{k,\Omega}$ and $[u]_{k,\Omega}$ denote norm
 135 and semi-norms for $H^k(\Omega) = W^{k,2}(\Omega)$. When there is no confusion, Ω may
 136 be dropped in the norm and semi-norms, e.g., $\|u\|_k := \|u\|_{k,\Omega}$.

- Inverse estimates for polynomials:

$$\|v_h\|_{k+1,e} \leq Ch^{-1} \|v_h\|_{k,e}, \quad \forall v_h \in V^h, k \geq 0.$$

- Elliptic regularity holds for the problem (3.1):

$$\|u\|_2 \leq C\|f\|_0$$

137

138

139

140

- Let Ω_h is a finite element mesh for Ω . For each element $e \in \Omega_h$, we denote $\bar{\mathbf{a}}_e = (\bar{a}_e^{ij})$ as an approximation to the average of \mathbf{a} on element e , i.e. $\bar{a}_e^{ij} = \frac{1}{\text{meas}(e)} \int_e a^{ij} d\mathbf{x}$. Specifically, we choose \bar{a}_e^{ij} as the the function value of a^{ij} at the center of element e . Then we define piece-wise constant function $\bar{\mathbf{a}}$ as

141

$$\bar{\mathbf{a}}(\mathbf{x}) = \bar{\mathbf{a}}_e, \quad \text{for } \mathbf{x} \in e.$$

142

- Define the projection operator $\hat{\Pi}_1 : \hat{u} \in L^1(\hat{K}) \rightarrow \hat{\Pi}_1 \hat{u} \in Q^1(\hat{K})$ by

143

$$(2.2) \quad \int_{\hat{K}} (\hat{\Pi}_1 \hat{u}) \hat{w} d\hat{\mathbf{x}} = \int_{\hat{K}} \hat{u} \hat{w} d\hat{\mathbf{x}}, \quad \forall \hat{w} \in Q^1(\hat{K}).$$

Observe that all degrees of freedom of $\hat{\Pi}_1 \hat{u}$ can be expressed as a linear combination of $\int_{\hat{K}} \hat{u} \hat{p} d\hat{\mathbf{x}}$ where $\hat{p}(\mathbf{x})$ takes the forms $1, \hat{x}_1, \hat{x}_2$, and $\hat{x}_1 \hat{x}_2$. This implies that the $H^1(\hat{K})$ (or $H^2(\hat{K})$) norm of $\hat{\Pi}_1 \hat{u}$ is dictated by $\int_{\hat{K}} \hat{u} \hat{p} d\hat{\mathbf{x}}$. Utilizing the Cauchy-Schwartz inequality, we deduce:

$$\left| \int_{\hat{K}} \hat{u} \hat{p} d\hat{\mathbf{x}} \right| \leq \|\hat{u}\|_{0,2,\hat{K}} \|\hat{p}\|_{0,2,\hat{K}} \leq C \|\hat{u}\|_{0,2,\hat{K}}$$

From which it follows that:

$$\|\Pi_1 \hat{u}\|_{1,2,\hat{K}} \leq C \|\hat{u}\|_{0,2,\hat{K}}$$

144

145

146

This establishes that $\hat{\Pi}_1$ acts as a continuous linear mapping from $L^2(\hat{K})$ to $H^1(\hat{K})$. Similarly, by extending this argument, we can also demonstrate that $\hat{\Pi}_1$ is a continuous linear mapping from $L^2(\hat{K})$ to $H^2(\hat{K})$.

- We denote all the the vertices of Ω_h inside Ω by $\mathbf{x}_j, j = 1, \dots, N_h$. We denote nodal basis functions in V_h by $\varphi_i, i = 1, \dots, N_h$, which are continuous in Ω , linear in each element e and

$$\varphi_i(\mathbf{x}_i) = 1, \quad \varphi_i(\mathbf{x}_j) = 0, \quad j \neq i.$$

147

148

149

150

2.2. Mixed quadrature. To analyze and impose the monotonicity of the stiffness matrix, we will use numerical quadrature rules to approximate integrals. As we will see, the choice of quadrature rules can significantly affect the monotonicity of the numerical schemes.

151

152

153

154

For a one-dimensional integral of function f over the interval $[0, 1]$, we can approximate $\int_0^1 f(\hat{x}) d\hat{x}$ using either the trapezoid rule, given by $\frac{f(0)+f(1)}{2}$, or the midpoint rule, $f(\frac{1}{2})$. Both quadrature offer second-order accuracy. We will use the linear combination of these two kinds of quadrature as follows:

155 (2.3)

$$\begin{aligned} \int_0^1 f(\hat{x}) d\hat{x} &\simeq \lambda \frac{f(0) + f(1)}{2} + (1 - \lambda) f\left(\frac{1}{2}\right) \\ &= \hat{\omega}_1 f(\hat{\xi}_1) + \hat{\omega}_2 f(\hat{\xi}_2) + \hat{\omega}_3 f(\hat{\xi}_1), \end{aligned}$$

156 where λ is a parameter to be determined and

$$157 \quad (2.4) \quad \hat{\omega}_1 = \frac{\lambda}{2}, \quad \hat{\omega}_2 = 1 - \lambda, \quad \hat{\omega}_3 = \frac{\lambda}{2}, \quad \hat{\xi}_1 = 0, \quad \hat{\xi}_2 = \frac{1}{2}, \quad \hat{\xi}_3 = 1.$$

158 When $\lambda = 1$, the mixed quadrature recovers the trapezoid rule and when $\lambda = 0$ the
159 mixed quadrature recovers the midpoint rule.

160 To approximate integration on square \hat{K} , we may use the mixed quadrature (2.3)
161 with different parameters λ^1 and λ^2 for different dimension x_1 and x_2 respectively.
162 By Fubini's theorem,

$$(2.5) \quad \int_{\hat{K}} f(\hat{\mathbf{x}}) d\hat{\mathbf{x}} = \int_0^1 \int_0^1 f(\hat{\mathbf{x}}) d\hat{\mathbf{x}} = \int_0^1 \left(\int_0^1 f(\hat{x}_1, \hat{x}_2) d\hat{x}_2 \right) d\hat{x}_1$$

$$163 \quad \simeq \int_0^1 \left(\sum_{q=1}^3 \hat{\omega}_q^2 f(\hat{x}_1, \hat{\xi}_q) \right) d\hat{x}_1 \simeq \sum_{p=1}^{r+1} \hat{\omega}_p^1 \left(\sum_{q=1}^{r+1} \hat{\omega}_q^2 f(\hat{\xi}_p, \hat{\xi}_q) \right) = \sum_{p=1}^3 \sum_{q=1}^3 \hat{\omega}_p^1 \hat{\omega}_q^2 f(\hat{\xi}_p, \hat{\xi}_q),$$

164 where ω_i^j are just ω_i while replacing λ with λ^j in (2.4) for $i = 1, 2, 3$, $j = 1, 2$.

165 On the reference element \hat{K} , for convenience, to denote the above quadrature
166 for integral approximation with parameter $\boldsymbol{\lambda} = (\lambda^1, \lambda^2)$, we will use the following
167 notation

$$168 \quad (2.6) \quad \int_{\hat{K}} \hat{f}(\hat{\mathbf{x}}) d_{\boldsymbol{\lambda}}^h \hat{\mathbf{x}} := \sum_{p=1}^3 \sum_{q=1}^3 \hat{\omega}_p^1 \hat{\omega}_q^2 f(\hat{\xi}_p, \hat{\xi}_q).$$

169 Given the quadrature parameter $\boldsymbol{\lambda}_e = (\lambda_e^1, \lambda_e^2)$, the quadrature approximation to
170 $\int_e f(\mathbf{x}) d\mathbf{x}$ is denoted as

$$171 \quad (2.7) \quad \int_e f(\mathbf{x}) d_{\boldsymbol{\lambda}_e}^h \mathbf{x} := \int_{\hat{K}} f \circ \mathbf{F}_e(\hat{\mathbf{x}}) d_{\boldsymbol{\lambda}_e}^h \hat{\mathbf{x}}.$$

172 Then we define the quadrature approximation over the entire domain Ω as

$$173 \quad (2.8) \quad \int_{\Omega} f d_{\boldsymbol{\lambda}_{\Omega}}^h \mathbf{x} := \sum_{e \in \Omega_h} \int_e f d_{\boldsymbol{\lambda}_e}^h \mathbf{x},$$

174 where $\boldsymbol{\lambda}_{\Omega} = (\boldsymbol{\lambda}_e)_{e \in \Omega_h}$ can be viewed as a vector-valued piece-wise constant function,
175 with values $\boldsymbol{\lambda}_e$ that differ across elements.

176 As a particular instance, $\int_{\Omega} f d_{\mathbf{1}}^h \mathbf{x}$ denote the case $\boldsymbol{\lambda}_e = (1, 0)$ for all $e \in \Omega_h$, i.e.
177 the integral on each element are approximated by the trapezoid rule in all directions.

178 **2.3. Quadrature error estimates.** The Bramble-Hilbert Lemma for Q^k poly-
179 nomials can be stated as follows, see Exercise 3.1 .1 and Theorem 4.1.3 in [4]:

180 **THEOREM 2.1.** *If a continuous linear mapping $\hat{\Pi} : H^{k+1}(\hat{K}) \rightarrow H^{k+1}(\hat{K})$ satis-
181 fies $\hat{\Pi}\hat{v} = \hat{v}$ for any $\hat{v} \in Q^k(\hat{K})$, then*

$$182 \quad (2.9) \quad \|\hat{u} - \hat{\Pi}\hat{u}\|_{k+1, \hat{K}} \leq C[\hat{u}]_{k+1, \hat{K}}, \quad \forall \hat{u} \in H^{k+1}(\hat{K}).$$

183 *Therefore if $l(\cdot)$ is a continuous linear form on the space $H^{k+1}(\hat{K})$ satisfying $l(\hat{v}) =$
184 $0, \forall \hat{v} \in Q^k(\hat{K})$, then*

$$185 \quad |l(\hat{u})| \leq C \|l\|'_{k+1, \hat{K}} [\hat{u}]_{k+1, \hat{K}}, \quad \forall \hat{u} \in H^{k+1}(\hat{K}),$$

186 *where $\|l\|'_{k+1, \hat{K}}$ is the norm in the dual space of $H^{k+1}(\hat{K})$.*

187 By applying Bramble-Hilbert Lemma, we have the following quadrature estimates.

188 LEMMA 2.2. For a sufficiently smooth function $a \in H^2(e)$, we have

$$189 \quad (2.10) \quad \int_e a d\mathbf{x} - \int_e a d^h \mathbf{x} = \mathcal{O}\left(h^{2+\frac{d}{2}}\right) [a]_{2,e} = \mathcal{O}\left(h^{2+d}\right) [a]_{2,\infty,e}$$

$$190 \quad (2.11) \quad \int_e a d\mathbf{x} - \int_e \bar{a}_e d\mathbf{x} = \mathcal{O}\left(h^{2+\frac{d}{2}}\right) [a]_{2,e} = \mathcal{O}\left(h^{2+d}\right) [a]_{2,\infty,e}$$

191

Proof. For any $\hat{f} \in H^2(\hat{K})$, since quadrature are represented by point values, with the Sobolev's embedding we have

$$|\hat{E}(\hat{f})| \leq C|\hat{f}|_{0,\infty,\hat{K}} \leq C\|\hat{f}\|_{2,2,\hat{K}}$$

Therefore $\hat{E}(\cdot)$ is a continuous linear form on $H^2(\hat{K})$ and $\hat{E}(\hat{f}) = 0$ if $\hat{f} \in Q^1(\hat{K})$. Then the Bramble-Hilbert lemma implies

$$|E(a)| = h^d |\hat{E}(\hat{a})| \leq Ch^d [\hat{a}]_{2,2,\hat{K}} = \mathcal{O}\left(h^{2+\frac{d}{2}}\right) [a]_{2,2,e} = \mathcal{O}\left(h^{2+d}\right) [a]_{2,\infty,e}$$

LEMMA 2.3. If $f \in H^2(\Omega)$, $\forall v_h \in V^h$, we have

$$(f, v_h) - \langle f, v_h \rangle_h = \mathcal{O}\left(h^2\right) \|f\|_2 \|v_h\|_1.$$

192 *Proof.* Applying Theorem 2.1, on element e , with $\frac{\partial^2 \hat{v}_h}{\partial^2 \hat{x}_i}$ vanish, we obtain:

$$\begin{aligned} E(fv) &= h^d \hat{E}(\hat{f}\hat{v}_h) \leq Ch^d [\hat{f}\hat{v}_h]_{2,2,\hat{K}} \\ &\leq Ch^d \left(|\hat{f}|_{2,2,\hat{K}} |\hat{v}_h|_{0,\infty,\hat{K}} + |\hat{f}|_{1,2,\hat{K}} |\hat{v}_h|_{1,\infty,\hat{K}} \right) \\ 193 \quad &\leq Ch^d \left(|\hat{f}|_{2,2,\hat{K}} |\hat{v}_h|_{0,2,\hat{K}} + |\hat{f}|_{1,2,\hat{K}} |\hat{v}_h|_{1,2,\hat{K}} \right) \\ &\leq Ch^2 (|f|_{2,2,e} |v_h|_{0,2,e} + |f|_{1,2,e} |v_h|_{1,2,e}) = \mathcal{O}\left(h^2\right) \|f\|_{2,e} \|v_h\|_{1,e}. \end{aligned}$$

By sum the above result over all elements of Ω_h , then we conclude with

$$(f, v_h) - \langle f, v_h \rangle_h = \mathcal{O}\left(h^2\right) \|f\|_2 \|v_h\|_1.$$

LEMMA 2.4. If $u \in H^3(e)$, for $i, j = 1, 2$, then $\forall v_h$,

$$\int_e u_{x_i}(v_h)_{x_j} d\mathbf{x} - \int u_{x_i}(v_h)_{x_j} d_{\lambda_e}^h \mathbf{x} = \mathcal{O}\left(h^2\right) \|u\|_{3,e} \|v_h\|_{2,e}.$$

194 *Proof.* Applying Theorem 2.1, we obtain:

$$\begin{aligned} E(u_{x_i}(v_h)_{x_j}) &= h^{d-2} \hat{E}(\hat{u}_{\hat{x}_i}(\hat{v}_h)_{\hat{x}_j}) \leq Ch^{d-2} [\hat{u}_{\hat{x}_i}(\hat{v}_h)_{\hat{x}_j}]_{2,2,\hat{K}} \\ &\leq Ch^{d-2} \left(|\hat{u}_{\hat{x}_i}|_{2,2,\hat{K}} |(\hat{v}_h)_{\hat{x}_j}|_{0,\infty,\hat{K}} + |\hat{u}_{\hat{x}_i}|_{1,2,\hat{K}} |(\hat{v}_h)_{\hat{x}_j}|_{1,\infty,\hat{K}} + |\hat{u}_{\hat{x}_i}|_{0,2,\hat{K}} |(\hat{v}_h)_{\hat{x}_j}|_{2,\infty,\hat{K}} \right) \\ 195 \quad &\leq Ch^{d-2} \left(|\hat{u}_{\hat{x}_i}|_{2,2,\hat{K}} |(\hat{v}_h)_{\hat{x}_j}|_{0,2,\hat{K}} + |\hat{u}_{\hat{x}_i}|_{1,2,\hat{K}} |(\hat{v}_h)_{\hat{x}_j}|_{1,2,\hat{K}} + |\hat{u}_{\hat{x}_i}|_{0,2,\hat{K}} |(\hat{v}_h)_{\hat{x}_j}|_{2,2,\hat{K}} \right) \\ &\leq Ch^{d-2} \left(|\hat{u}|_{3,2,\hat{K}} |\hat{v}_h|_{1,2,\hat{K}} + |\hat{u}|_{2,2,\hat{K}} |\hat{v}_h|_{2,2,\hat{K}} \right). \end{aligned}$$

196 where the second last inequality is implied by the equivalence of norms over $Q^1(\hat{K})$
197 and in the last inequality we use the fact that the third derivative of Q^1 polynomial
198 vanish.

199 Therefore,

$$200 \quad E(u_{x_i}(v_h)_{x_j}) \leq Ch^2 (|u|_{3,2,e}|v_h|_{1,2,e} + |u|_{2,2,e}|v_h|_{2,2,e}) = \mathcal{O}(h^2) \|u\|_{3,e} \|v_h\|_{2,e}. \quad \square$$

LEMMA 2.5. *If $f \in H^2(\Omega)$ or $f \in V^h$, $\forall v_h$, we have*

$$(f, v_h) - \langle f, v_h \rangle_h = \mathcal{O}(h) \|f\|_2 \|v_h\|_0.$$

201 *Proof.* As in the proof of Lemma 2.3, we have

$$202 \quad E(fv) = \mathcal{O}(h^2) \|f\|_{2,e} \|v_h\|_{1,e}. \quad \square$$

By applying the inverse estimate to polynomial v_h , we have

$$E(fv) = \mathcal{O}(h) \|f\|_{2,e} \|v_h\|_{0,e}.$$

Summing the previous result across all elements in Ω_h , we conclude:

$$(f, v_h) - \langle f, v_h \rangle_h = \mathcal{O}(h) \|f\|_2 \|v_h\|_0.$$

203 **3. The Q^1 finite element method and its monotonicity.** In this section, we
204 give a derivation of the Q^1 finite element scheme and then discuss its monotonicity.

205 **3.1. Derivation of the scheme.** The variational form of (1.1) is to find $u \in$
206 $H_0^1(\Omega)$ satisfying

$$207 \quad (3.1) \quad \mathcal{A}(u, v) = (f, v), \quad \forall v \in H_0^1(\Omega),$$

208 where $\mathcal{A}(u, v) = \int_{\Omega} \mathbf{a} \nabla u \cdot \nabla v d\mathbf{x} + \int_{\Omega} cuv d\mathbf{x}$, $(f, v) = \int_{\Omega} f v d\mathbf{x}$.

Let $V_0^h \subseteq H_0^1(\Omega)$ be the continuous finite element space consisting of piece-wise Q^1 polynomials. To have a second-order monotone method, we first approximate the matrix coefficients $\mathbf{a} = (a^{ij}(\mathbf{x}))$ by either its average $\frac{1}{\text{meas}(e)} \int_e \mathbf{a} d\mathbf{x}$ or its middle point value on each element e . The approximation is denoted by $\bar{\mathbf{a}}_e$. Then we get the modified bilinear form

$$\bar{\mathcal{A}}(u, v) = \int_{\Omega} \bar{\mathbf{a}} \nabla u \cdot \nabla v d\mathbf{x} + \int_{\Omega} cuv d\mathbf{x},$$

209 where $\bar{\mathbf{a}} = (\bar{\mathbf{a}}_e)_{e \in \Omega_h}$. In practice, we take $\bar{\mathbf{a}}_e$ to be the middle point value of $\bar{\mathbf{a}}$ on
210 element e for smooth enough \mathbf{a} and fine enough mesh.

211 By approximating integrals in $\bar{\mathcal{A}}(u_h, v_h)$ with quadrature specified in (2.8), along
212 with designated quadrature parameter λ_{Ω} , we derive the following numerical scheme:
213 find $u_h \in V_0^h$ satisfying

$$214 \quad (3.2) \quad \mathcal{A}_h(u_h, v_h) = \langle f, v_h \rangle_h, \quad \forall v_h \in V_0^h,$$

215 where the approximated bilinear form is defined as

$$216 \quad (3.3) \quad \mathcal{A}_h(u_h, v_h) := \int_{\Omega} \bar{\mathbf{a}} \nabla u_h \cdot \nabla v_h d_{\lambda}^h \mathbf{x} + \int_{\Omega} cu_h v_h d_1^h \mathbf{x}$$

217 and the right hand side is

$$218 \quad (3.4) \quad \langle f, v_h \rangle_h := \int_{\Omega} f v_h d_1^h \mathbf{x}.$$

219 Of course, the quadrature parameter $\lambda = (\lambda^1, \lambda^2)$ on each element need to be
220 determined for the quadrature (2.7).

221 It is not obvious that the numerical solution u_h is an accurate approximation of
222 the exact solution u as $\bar{\mathbf{a}}$ varies depending on the mesh.

223 **3.2. Monotonicity.** Let $A = (\mathcal{A}_h(\nabla\varphi_i, \nabla\varphi_j))$ be the stiffness matrix of our Q^1
 224 scheme (3.2) for equation (1.1). To have the monotonicity, we enforce the stiffness
 225 matrix A to be a M -matrix. We are interested in conditions for A to be an M -matrix.
 226 Recall a sufficient condition for M -matrix, see condition C_{10} in [29]:

227 **LEMMA 3.1.** *For a real irreducible square matrix A with positive diagonal entries*
 228 *and non-positive off-diagonal entries, A is a nonsingular M -matrix if all the row sums*
 229 *of A are non-negative and at least one row sum is positive.*

230 Then we have the following result on the uniform rectangular mesh.

231 **THEOREM 3.2.** *Assume $\forall e \in \Omega_h$, $|\bar{a}_e^{12}| \leq \min\{\bar{a}_e^{11}, \bar{a}_e^{22}\}$. Then for the Q^1 scheme*
 232 *given by (3.2) for the elliptic equation (1.1) on uniform rectangular mesh, the stiffness*
 233 *matrix is a M -matrix, provided the quadrature parameters for each element e are*
 234 *chosen as:*

$$235 \quad (3.5) \quad \lambda_e^1, \lambda_e^2 \in \left(\frac{|\bar{a}_e^{11} - \bar{a}_e^{22}|}{\bar{a}_e^{11} + \bar{a}_e^{22}}, 1 - \frac{2|\bar{a}_e^{12}|}{\bar{a}_e^{11} + \bar{a}_e^{22}} \right].$$

236 When $|\bar{a}_e^{12}| = \min\{\bar{a}_e^{11}, \bar{a}_e^{22}\}$, (3.5) means we take λ_e^1, λ_e^2 to be the upper bound of the
 237 interval, i.e. $1 - \frac{2|\bar{a}_e^{12}|}{\bar{a}_e^{11} + \bar{a}_e^{22}}$.

238 *Proof.* First, we consider the following quadrature approximation results on the
 239 reference element \hat{K} . With quadrature (2.6) and quadrature parameter $\lambda_e = (\lambda_e^1, \lambda_e^2)$,
 240 we have

$$241 \quad \langle \bar{\mathbf{a}}\nabla\phi_{0,0}, \nabla\phi_{0,1} \rangle_h = \langle \bar{\mathbf{a}}\nabla\phi_{1,1}, \nabla\phi_{1,0} \rangle_h = -\frac{1}{4}(\lambda_e^2\bar{a}_e^{11} + \lambda_e^1\bar{a}_e^{22}) + \frac{1}{4}(\bar{a}_e^{11} - \bar{a}_e^{22}),$$

$$242 \quad \langle \bar{\mathbf{a}}\nabla\phi_{0,0}, \nabla\phi_{1,0} \rangle_h = \langle \bar{\mathbf{a}}\nabla\phi_{0,1}, \nabla\phi_{1,1} \rangle_h = -\frac{1}{4}(\lambda_e^2\bar{a}_e^{11} + \lambda_e^1\bar{a}_e^{22}) + \frac{1}{4}(\bar{a}_e^{22} - \bar{a}_e^{11}),$$

$$243 \quad \langle \bar{\mathbf{a}}\nabla\phi_{0,0}, \nabla\phi_{1,1} \rangle_h = -\frac{1}{4}((1 - \lambda_e^2)\bar{a}_e^{11} + (1 - \lambda_e^1)\bar{a}_e^{22}) - \frac{1}{2}\bar{a}_e^{12},$$

$$244 \quad \langle \bar{\mathbf{a}}\nabla\phi_{0,1}, \nabla\phi_{1,0} \rangle_h = -\frac{1}{4}((1 - \lambda_e^2)\bar{a}_e^{11} + (1 - \lambda_e^1)\bar{a}_e^{22}) + \frac{1}{2}\bar{a}_e^{12}.$$

246 With (3.5) and the assumption $|\bar{a}_e^{12}| \leq \min\{\bar{a}_e^{11}, \bar{a}_e^{22}\}$, we have
 (3.6)

$$\langle \bar{\mathbf{a}}\nabla\phi_{0,0}, \nabla\phi_{0,1} \rangle_h = \langle \bar{\mathbf{a}}\nabla\phi_{1,1}, \nabla\phi_{1,0} \rangle_h \in \left[\frac{1}{2}(|\bar{a}_e^{12}| - \bar{a}_e^{22}), \frac{1}{4}(\bar{a}_e^{11} - \bar{a}_e^{22} - |\bar{a}_e^{11} - \bar{a}_e^{22}|) \right],$$

$$\langle \bar{\mathbf{a}}\nabla\phi_{0,0}, \nabla\phi_{1,0} \rangle_h = \langle \bar{\mathbf{a}}\nabla\phi_{0,1}, \nabla\phi_{1,1} \rangle_h \in \left[\frac{1}{2}(|\bar{a}_e^{12}| - \bar{a}_e^{11}), \frac{1}{4}(\bar{a}_e^{22} - \bar{a}_e^{11} - |\bar{a}_e^{11} - \bar{a}_e^{22}|) \right],$$

$$247 \quad \langle \bar{\mathbf{a}}\nabla\phi_{0,0}, \nabla\phi_{1,1} \rangle_h \in \left(-\frac{1}{2}(\min\{\bar{a}_e^{11}, \bar{a}_e^{22}\} - \bar{a}_e^{12}), -\frac{1}{2}(|\bar{a}_e^{12}| + \bar{a}_e^{12}) \right],$$

$$\langle \bar{\mathbf{a}}\nabla\phi_{0,1}, \nabla\phi_{1,0} \rangle_h \in \left(-\frac{1}{2}(\min\{\bar{a}_e^{11}, \bar{a}_e^{22}\} + \bar{a}_e^{12}), -\frac{1}{2}(|\bar{a}_e^{12}| - \bar{a}_e^{12}) \right],$$

248 which are all non-positive. Again, when $|\bar{a}_e^{12}| = \min\{\bar{a}_e^{11}, \bar{a}_e^{22}\}$, we will take the above
 249 values as the bound of the closed side of the interval.

250 Given $j \in \{1, \dots, N_h\}$, consider the corresponding node x_j . Obviously, if both x_i

251 and x_j are vertices of the same elements e ,

$$\begin{aligned}
A_{ij} &= \mathcal{A}_h(\varphi_j, \varphi_i) \\
&= \sum_{e \in \Omega_h} \int_e \bar{\mathbf{a}} \nabla \varphi_j \cdot \nabla \varphi_i d_{\lambda_e}^h \mathbf{x} + \int_e c \varphi_j \varphi_i d_1^h \mathbf{x} \\
252 \quad (3.7) \quad &= \sum_{e \in \Omega_h} \int_{\hat{K}} \bar{\mathbf{a}} \hat{\nabla} \hat{\varphi}_j \cdot \hat{\nabla} \hat{\varphi}_i d_{\lambda_e}^h \hat{\mathbf{x}} + \int_{\hat{K}} \hat{c} \hat{\varphi}_j \hat{\varphi}_i d_1^h \hat{\mathbf{x}} \\
&= \sum_{i,j \in e} \int_{\hat{K}} \bar{\mathbf{a}} \hat{\nabla} \hat{\varphi}_j \cdot \hat{\nabla} \hat{\varphi}_i d_{\lambda_e}^h \hat{\mathbf{x}} + \int_{\hat{K}} \hat{c} \hat{\varphi}_j \hat{\varphi}_i d_1^h \hat{\mathbf{x}}
\end{aligned}$$

253 where $\sum_{i,j \in e}$ means summation over all elements e containing both vertices i and j .

254 Notice that $\int_{\hat{K}} \hat{c} \hat{\varphi}_j \hat{\varphi}_i d_1^h \hat{\mathbf{x}}$ vanish if $i \neq j$ and $\int_{\hat{K}} \bar{\mathbf{a}} \hat{\nabla} \hat{\varphi}_j \cdot \hat{\nabla} \hat{\varphi}_i d_{\lambda_e}^h \hat{\mathbf{x}}$ aligns with one
255 of the values in (3.6) depending on their relative positions. Therefore, for $i \neq j$, with
256 (3.5) and the assumption $|\bar{a}_e^{12}| \leq \min\{\bar{a}_e^{11}, \bar{a}_e^{22}\}$ we have

$$257 \quad (3.8) \quad A_{ij} = \sum_{i,j \in e} \int_{\hat{K}} \bar{\mathbf{a}} \hat{\nabla} \hat{\varphi}_j \cdot \hat{\nabla} \hat{\varphi}_i d_{\lambda_e}^h \hat{\mathbf{x}} \leq 0.$$

If \mathbf{x}_i has no neighboring node on the boundary, then the i -th row sum of A is non-negative:

$$\sum_j A_{ij} = \sum_{j=0}^{N_h} \mathcal{A}_h(\varphi_j, \varphi_i) = \mathcal{A}_h(1, \varphi_i) = C c_i \geq 0,$$

258 where C is a certain positive number and $c_i = c(\mathbf{x}_i) \geq 0$. Therefore, $A_{ii} \geq \sum_{j \neq i} |A_{ij}|$.

259 When \mathbf{x}_i has a neighboring node on the boundary, we do have $A_{ii} \geq \sum_{j \neq i} |A_{ij}|$.
260 When \mathbf{x}_i has two neighboring node on the boundary, based on (3.6), in the stencil of
261 x_i , one of the corresponding coefficients of the two neighboring nodes on the boundary
262 must be negative, and it is not in $A_{i,\cdot}$, then $\sum_j A_{ij} > 0$, i.e. $A_{ii} > \sum_{j \neq i} |A_{ij}|$.

263 Therefore, we conclude the proof. \square

264 **REMARK 1.** For each element e , the choice in (3.5) make $\lambda_e^1, \lambda_e^2 > 0$, which im-
265 plies the V^h -ellipticity of the bilinear form (3.3) discussed in Section 4.2. Therefore,
266 we can assure of V^h -ellipticity and the stiffness matrix being an M -matrix simultane-
267 ously.

268 **REMARK 2.** The constraint on the coefficient, $|\bar{a}_e^{12}| \leq \min\{\bar{a}_e^{11}, \bar{a}_e^{22}\}$, aligns with
269 the condition for rendering the stiffness matrix as an M -matrix in the seven-point
270 stencil control volume method with optimal optimal monotonicity region in the case
271 of homogeneous medium and uniform mesh in [28]. In [27], the authors show that
272 a three-by-three stencil can be used to construct monotone finite difference schemes
273 under the assumption $|a^{12}| < \min\{a^{11}, a^{22}\}$.

274 **4. Convergence of the Q^1 finite element method with mixed quadrature.**
275 In this section, we prove the second-order accuracy of the scheme (3.2) on
276 uniform rectangular mesh. For convenience, in this section, we may drop the sub-
277 script h in a test function $v_h \in V^h$. When there is no confusion, we may also drop dx
278 or $d\hat{\mathbf{x}}$ in an integral.

279 **4.1. Approximation error estimate of bilinear forms.** In this subsection,
280 we estimate the approximation error of $\mathcal{A}_h(u, v)$ to $\mathcal{A}(u, v)$.

281 THEOREM 4.1. Assume $a^{ij}, c \in W^{2,\infty}(\Omega)$ for $i, j = 1, 2$ and $u \in H^3(\Omega)$, then
 282 $\forall v \in V^h$, on element e , we have

$$283 \quad (4.1) \quad \int_e (\mathbf{a} \nabla u) \cdot \nabla v d\mathbf{x} - \int_e (\bar{\mathbf{a}}_e \nabla u) \cdot \nabla v d_{\lambda_e}^h \mathbf{x} = \mathcal{O}(h^2) \|u\|_{3,e} \|v\|_{2,e},$$

$$284 \quad (4.2) \quad \int_e c u v d\mathbf{x} - \int_e c u v d_1^h \mathbf{x} = \mathcal{O}(h^2) \|u\|_{2,e} \|v\|_{2,e}.$$

285
 286 *Proof.* For $k, l = 1, 2$ and function $a \in W^{2,\infty}(e)$, we have

$$287 \quad (4.3) \quad \begin{aligned} & \int_e a u_{x_k} v_{x_l} d\mathbf{x} - \int_e \bar{a}_e u_{x_k} v_{x_l} d_{\lambda_e}^h \mathbf{x} \\ &= \int_e (a - \bar{a}_e) u_{x_k} v_{x_l} d\mathbf{x} + \bar{a}_e \left(\int_e u_{x_k} v_{x_l} d\mathbf{x} - \int_e u_{x_k} v_{x_l} d_{\lambda_e}^h \mathbf{x} \right) \\ &= \int_e (a - \bar{a}_e) u_{x_k} v_{x_l} d\mathbf{x} + \bar{a}_e E(u_{x_k} v_{x_l}). \end{aligned}$$

288 For the first term,

$$289 \quad (4.4) \quad \begin{aligned} & \int_e (a - \bar{a}_e) u_{x_k} v_{x_l} d\mathbf{x} \\ &= \int_e (a - \bar{a}_e) (u_{x_k} v_{x_l} - \overline{u_{x_k} v_{x_l}}) d\mathbf{x} + \int_e (a - \bar{a}_e) \overline{u_{x_k} v_{x_l}} d\mathbf{x} \\ &\leq \|a - \bar{a}_e\|_{0,\infty,e} \|u_{x_k} v_{x_l} - \overline{u_{x_k} v_{x_l}}\|_{0,1,e} + \frac{1}{\text{meas}(e)} \int_e (a - \bar{a}_e) d\mathbf{x} \int_e u_{x_k} v_{x_l} d\mathbf{x}. \end{aligned}$$

290 By Poincare inequality and Cauchy-Schwartz inequality, we have

$$291 \quad (4.5) \quad \begin{aligned} & \|a - \bar{a}_e\|_{0,\infty,e} \|u_{x_k} v_{x_l} - \overline{u_{x_k} v_{x_l}}\|_{0,1,e} \\ &= \mathcal{O}(h^2) \|a\|_{1,\infty,e} \|\nabla(u_{x_k} v_{x_l})\|_{0,1,e} = \mathcal{O}(h^2) \|u\|_{2,e} \|v\|_{2,e}. \end{aligned}$$

292 By Lemma 2.2 and Cauchy-Schwartz inequality

$$293 \quad (4.6) \quad \begin{aligned} & \frac{1}{\text{meas}(e)} \int_e (a - \bar{a}_e) d\mathbf{x} \int_e u_{x_k} v_{x_l} d\mathbf{x} \\ &= \frac{h^{2+d}}{\text{meas}(e)} [a]_{2,\infty,e} \|u_{x_k}\|_{0,e} \|v_{x_l}\|_{0,e} = \mathcal{O}(h^2) \|u\|_{1,e} \|v\|_{1,e} \end{aligned}$$

294 where in the last equation $\text{meas}(e) = \mathcal{O}(h^d)$ is also used. Therefore, we have the
 295 estimate of the first term of (4.3):

$$296 \quad (4.7) \quad \int_e (a - \bar{a}_e) u_{x_k} v_{x_l} d\mathbf{x} = \mathcal{O}(h^2) \|a\|_{2,\infty,e} \|u\|_{2,e} \|v\|_{2,e}.$$

297 For the second term of (4.3), by Lemma 2.4, we obtain

$$298 \quad (4.8) \quad \int_e \bar{a}_e u_{x_k} v_{x_l} d\mathbf{x} - \int_e \bar{a}_e u_{x_l} v_{x_l} d_{\lambda_e}^h \mathbf{x} = \mathcal{O}(h^2) \|a\|_{0,\infty,e} \|u\|_{3,e} \|v\|_{2,e},$$

299 which together with (4.7) imply the estimate of (4.3):

$$300 \quad (4.9) \quad \int_e a u_{x_k} v_{x_l} d\mathbf{x} - \int_e \bar{a}_e u_{x_k} v_{x_l} d_{\lambda_e}^h \mathbf{x} = \mathcal{O}(h^2) \|a\|_{2,\infty,e} \|u\|_{3,e} \|v\|_{2,e}.$$

301 Therefore, we have

$$302 \quad (4.10) \quad \int_e (\mathbf{a}(\mathbf{x}) \nabla u) \cdot \nabla v d\mathbf{x} - \int_e (\bar{\mathbf{a}}(\mathbf{x}) \cdot \nabla u) \nabla v d_{\lambda_e}^h \mathbf{x} = \mathcal{O}(h^2) \|\mathbf{a}\|_{2,\infty,e} \|u\|_{3,e} \|v\|_{2,e}.$$

303 Similarly we have

$$304 \quad (4.11) \quad \int_e c v d\mathbf{x} - \int_e c v d_1^h \mathbf{x} = \mathcal{O}(h^2) \|c\|_{2,\infty,e} \|u\|_{2,e} \|v\|_{2,e}. \quad \square$$

305 We also have

LEMMA 4.2. *Assume $a^{ij}, c \in W^{2,\infty}(\Omega)$ for $i, j = 1, 2$. We have*

$$A(v_h, w_h) - A_h(v_h, w_h) = \mathcal{O}(h) \|v_h\|_2 \|w_h\|_1, \quad \forall v_h, w_h \in V^h$$

306 *Proof.* By Theorem 4.1 and noticing that the third derivative of Q^1 polynomial
307 vanish, we have

$$308 \quad (4.12) \quad \int_e (\mathbf{a} \nabla v_h) \cdot \nabla w_h d\mathbf{x} - \int_e (\bar{\mathbf{a}}_e \nabla v_h) \cdot \nabla w_h d_{\lambda_e}^h \mathbf{x} = \mathcal{O}(h^2) \|v_h\|_{2,e} \|w_h\|_{2,e},$$

$$309 \quad (4.13) \quad \int_e c v_h w_h d\mathbf{x} - \int_e c v_h w_h d_1^h \mathbf{x} = \mathcal{O}(h^2) \|v_h\|_{2,e} \|w_h\|_{2,e}.$$

310

311 By applying the inverse estimate to polynomial z_h , we get

$$312 \quad (4.14) \quad \int_e (\mathbf{a} \nabla v_h) \cdot \nabla w_h d\mathbf{x} - \int_e (\bar{\mathbf{a}}_e \nabla v_h) \cdot \nabla w_h d_{\lambda_e}^h \mathbf{x} = \mathcal{O}(h) \|v_h\|_{2,e} \|w_h\|_{1,e},$$

$$313 \quad (4.15) \quad \int_e c v_h w_h d\mathbf{x} - \int_e c v_h w_h d_1^h \mathbf{x} = \mathcal{O}(h) \|v_h\|_{2,e} \|w_h\|_{1,e}.$$

314

315 Then by summing over all the elements we get prove the Lemma. \square

316 **4.2. V^h -ellipticity and the dual problem.** In order to prove the convergence
317 results of the scheme (3.2), we need A_h satisfies V^h -ellipticity:

$$318 \quad (4.16) \quad \forall v_h \in V_0^h, \quad C \|v_h\|_1^2 \leq A_h(v_h, v_h).$$

319 By following the proof of Lemma 5.1 in [15], we have

LEMMA 4.3. *Assume the eigenvalues of \mathbf{a} have a uniform positive lower bound and a uniform upper bound and c have a upper bound. If there exists lower bound $\lambda_0 > 0$ such that $\forall e \in \Omega_h$, the quadrature parameter $\lambda_e^1, \lambda_e^2 > \lambda_0$, then there are two constants $C_1, C_2 > 0$ independent of mesh size h such that*

$$\forall v_h \in V_0^h, \quad C_1 \|v_h\|_1^2 \leq A_h(v_h, v_h) \leq C_2 \|v_h\|_1^2.$$

Proof. For element e , at first we map all the functions to the reference element \hat{K} . Let $Z_{0,\hat{K}}$ denote the set of vertices on the reference element \hat{K} . We notice that the set $Z_{0,\hat{K}}$ is a $Q^1(\hat{K})$ -unisolvent subset. Since the weights of trapezoid rule are strictly positive, we have

$$\forall \hat{p} \in Q^1(\hat{K}), \quad \sum_{i=1}^2 \int_{\hat{K}} \hat{p}_{\hat{x}^i}^2 d_1^h \hat{\mathbf{x}} = 0 \implies \hat{p}_{\hat{x}^i} = 0 \text{ at } Z_{0,\hat{K}},$$

where $i = 1, 2$. As a consequence, $\sum_{i=1}^2 \int_{\hat{K}} \hat{p}_{\hat{x}_i}^2 d_1^h \hat{\mathbf{x}}$ defines a norm over the quotient space $Q^1(\hat{K})/Q^0(\hat{K})$. Since that $|\cdot|_{1,\hat{K}}$ is also a norm over the same quotient space, by the equivalence of norms over a finite dimensional space, we have

$$\forall \hat{p} \in Q^1(\hat{K}), \quad C_1 |\hat{p}|_{1,\hat{K}}^2 \leq \sum_{i=1}^2 \int_{\hat{K}} \hat{p}_{\hat{x}_i}^2 d_1^h \hat{\mathbf{x}} \leq C_2 |\hat{p}|_{1,\hat{K}}^2$$

As the quadrature parameter $\lambda_e^1, \lambda_e^2 \geq \lambda_0 \geq 0$, we have

$$C_1 |\hat{v}_h|_{1,\hat{K}}^2 \leq C_1 \sum_{i=1}^2 \int_{\hat{K}} (\hat{v}_h)_{\hat{x}_i}^2 d_1^h \hat{\mathbf{x}} \leq \int_{\hat{K}} (\bar{\mathbf{a}}_e^{ij} \nabla \hat{v}_h) \cdot \nabla \hat{v}_h d_{\lambda_e}^h \hat{\mathbf{x}} + \int_{\hat{K}} \hat{c} \hat{v}_h^2 d_1^h \hat{\mathbf{x}} \leq C_2 \|\hat{v}_h\|_{1,\hat{K}}^2.$$

320 Mapping these back to the original element e and summing over all elements, by
 321 the equivalence of two norms $|\cdot|_1$ and $\|\cdot\|_1$ for the space $H_0^1(\Omega) \supset V_0^h$, we get the
 322 conclusion. \square

323 In the following part, we assume the assumption of Lemma 4.3 is fulfilled, i.e. the
 324 V^h -ellipticity holds.

325 In order to apply the Aubin-Nitsche duality argument for establishing convergence
 326 of function values, we need certain estimates on a proper dual problem.

327 Define $\theta := u - u_h$ and consider the dual problem: find $w \in H_0^1(\Omega)$ satisfying

$$328 \quad (4.17) \quad A^*(w, v) = (\theta, v), \quad \forall v \in H_0^1(\Omega),$$

where $A^*(\cdot, \cdot)$ is the adjoint bilinear form of $A(\cdot, \cdot)$ such that

$$A^*(u, v) = A(v, u) = (\mathbf{a} \nabla v, \nabla u) + (cv, u).$$

329 Although here the bilinear form we considered is symmetric i.e. $A(\cdot, \cdot) = A^*(\cdot, \cdot)$, we
 330 still use $A^*(\cdot, \cdot)$ for abstractness.

331 Let $w_h \in V_0^h$ be the solution to

$$332 \quad (4.18) \quad A_h^*(w_h, v_h) = (\theta, v_h), \quad \forall v_h \in V_0^h.$$

333 Notice that the right hand side of (4.18) is different from the right hand side of
 334 the scheme (3.2).

335 We have the following standard estimates on w_h for the dual problem.

336 LEMMA 4.4. Assume $a^{ij}, c \in W^{2,\infty}(\Omega)$ and $u \in H^3(\Omega), f \in H^2(\Omega)$. Let w be
 337 defined in (4.17), w_h be defined in (4.18). With elliptic regularity and V^h -ellipticity
 338 hold, we have

$$339 \quad (4.19) \quad \begin{aligned} \|w - w_h\|_1 &\leq Ch \|w\|_2 \\ \|w_h\|_2 &\leq C \|\theta\|_0. \end{aligned}$$

Proof. By V^h -ellipticity, we have $C_1 \|w_h - v_h\|_1^2 \leq A_h^*(w_h - v_h, w_h - v_h)$. By the definition of the dual problem (4.17), we have

$$A_h^*(w_h, w_h - v_h) = (\theta, w_h - v_h) = A^*(w, w_h - v_h), \quad \forall v_h \in V_0^h.$$

Therefore $\forall v_h \in V_0^h$, by Lemma 4.2, we have

$$\begin{aligned} &C_1 \|w_h - v_h\|_1^2 \leq A_h^*(w_h - v_h, w_h - v_h) \\ &= A^*(w - v_h, w_h - v_h) + [A_h^*(w_h, w_h - v_h) - A^*(w, w_h - v_h)] + [A^*(v_h, w_h - v_h) - A_h^*(v_h, w_h - v_h)] \\ &= A^*(w - v_h, w_h - v_h) + [A(w_h - v_h, v_h) - A_h(w_h - v_h, v_h)] \\ &\leq C \|w - v_h\|_1 \|w_h - v_h\|_1 + Ch \|v_h\|_2 \|w_h - v_h\|_1, \end{aligned}$$

340 which implies

$$341 \quad (4.20) \quad \|w - w_h\|_1 \leq \|w - v_h\|_1 + \|w_h - v_h\|_1 \leq C \|w - v_h\|_1 + Ch \|v_h\|_2.$$

342 Now consider $\Pi_1 w \in V_0^h$ where Π_1 is the piece-wise Q^1 projection and its defini-
343 tion on each element is defined through (2.2) on the reference element. By Theorem
344 2.1 on the projection error, we have

$$345 \quad (4.21) \quad \|w - \Pi_1 w\|_1 \leq Ch \|w\|_2, \quad \|w - \Pi_1 w\|_2 \leq C \|w\|_2,$$

346 which implies

$$347 \quad (4.22) \quad \|\Pi_1 w\|_2 \leq \|w\|_2 + \|w - \Pi_1 w\|_2 \leq C \|w\|_2.$$

348 By setting $v_h = \Pi_1 w$, using (4.20), (4.21) and (4.22), we have

$$349 \quad (4.23) \quad \|w - w_h\|_1 \leq C \|w - \Pi_1 w\|_1 + Ch \|\Pi_1 w\|_2 \leq Ch \|w\|_2.$$

350 By (4.21) and (4.23), we also have

$$351 \quad (4.24) \quad \|w_h - \Pi_1 w\|_1 \leq \|w - \Pi_1 w\|_1 + \|w - w_h\|_1 \leq Ch \|w\|_2.$$

352 By the inverse estimate on the piece-wise polynomial $w_h - \Pi_1 w$, we get

$$353 \quad (4.25) \quad \|w_h\|_2 \leq \|w_h - \Pi_1 w\|_2 + \|\Pi_1 w - w\|_2 + \|w\|_2 \leq Ch^{-1} \|w_h - \Pi_1 w\|_1 + C \|w\|_2. \square$$

With (4.24), (4.25) and the elliptic regularity $\|w\|_2 \leq C \|\theta\|_0$, we get

$$\|w_h\|_2 \leq C \|w\|_2 \leq C \|\theta\|_0.$$

354 **4.3. Convergence results.** In this section, we initially establish the error es-
355 timate for $\|u - u_h\|_{1,\Omega}$. Subsequently, we demonstrate that the Q^1 finite element
356 method, as given by (3.2), achieves second-order accuracy for function values.

357 We have the estimate of the error $\|u - u_h\|_{1,\Omega}$ as follows:

THEOREM 4.5. *Assume $a^{ij}, c \in W^{2,\infty}(\Omega)$ and $u \in H^2(\Omega), f \in H^2(\Omega)$. With elliptic regularity and V^h -ellipticity hold, we have*

$$\|u - u_h\|_{1,\Omega} = \mathcal{O}(h) (\|u\|_{2,\Omega} + \|f\|_{2,\Omega}).$$

358 *Proof.* By the First Strang Lemma,

$$359 \quad (4.26) \quad \|u - u_h\|_{1,\Omega} \leq C \left(\inf_{v_h \in V^h} \left\{ \|u - v_h\|_{1,\Omega} + \sup_{w_h \in V^h} \frac{|\mathcal{A}(v_h, w_h) - \mathcal{A}_h(v_h, w_h)|}{\|w_h\|_{1,\Omega}} \right\} + \right. \\ \left. + \sup_{w_h \in V^h} \frac{|\langle f, w_h \rangle_h - (f, w_h)|}{\|w_h\|_{1,\Omega}} \right).$$

360 By Lemma 4.2, we have:

$$361 \quad \frac{|\mathcal{A}(v_h, w_h) - \mathcal{A}_h(v_h, w_h)|}{\|w_h\|_{1,\Omega}} = \frac{\mathcal{O}(h) \|v_h\|_{2,\Omega} \|w_h\|_{1,\Omega}}{\|w_h\|_{1,\Omega}} = \mathcal{O}(h) \|v_h\|_{2,\Omega}.$$

By Lemma 2.3, we have

$$\sup_{w_h \in V^h} \frac{|\langle f, w_h \rangle_h - (f, w_h)|}{\|w_h\|_{1,\Omega}} = \frac{\mathcal{O}(h^2)\|f\|_{2,\Omega}\|w_h\|_{1,\Omega}}{\|w_h\|_{1,\Omega}} = \mathcal{O}(h^2)\|f\|_{2,\Omega}.$$

By the approximation property of piece-wise Q^1 polynomials, □

$$\|u - u_h\|_{1,\Omega} = \mathcal{O}(h)(\|u\|_{2,\Omega} + \|f\|_{2,\Omega}).$$

In the following part we prove the Aubin-Nitsche Lemma up to the quadrature error for establishing convergence of function values.

THEOREM 4.6. *Assume $a^{ij}, c \in W^{2,\infty}(\Omega)$ and $u(\mathbf{x}) \in H^3(\Omega), f \in H^2(\Omega)$. Assume V^h ellipticity holds. Then the numerical solution from scheme (3.2) u_h is a 2-th order accurate approximation to the exact solution u :*

$$\|u_h - u\|_{0,\Omega} = \mathcal{O}(h^2)(\|u\|_{2,\Omega} + \|f\|_{2,\Omega}).$$

Proof. With $\theta = u - u_h \in H_0^1(\Omega)$, we have

$$(4.27) \quad \|\theta\|_0^2 = (\theta, \theta) = A(\theta, w) = A(u - u_h, w_h) + A(u - u_h, w - w_h)$$

For the first term (4.27), by Lemma 4.1, we have

$$(4.28) \quad \begin{aligned} A(u - u_h, w_h) &= [A(u, w_h) - A_h(u_h, w_h)] + [A_h(u_h, w_h) - A(u_h, w_h)] \\ &= (f, w_h) - \langle f, w_h \rangle_h + \mathcal{O}(h^2)\|u_h\|_3\|w_h\|_2 \\ &= \mathcal{O}(h^2)\|f\|_2\|w_h\|_1 + \mathcal{O}(h^2)\|u_h\|_2\|w_h\|_2 \\ &= \mathcal{O}(h^2)(\|f\|_2 + \|u_h\|_2)\|\theta\|_0, \end{aligned}$$

where in the second last equation Lemma 2.3 and the fact the third derivative of Q^1 polynomials vanish are used. As the estimate of $\|w_h\|_2$ and $\|w\|_2$ in the proof of Lemma 4.4, we have

$$(4.29) \quad \begin{aligned} \|u_h\|_2 &\leq \|u_h - \Pi_1 u\|_2 + \|\Pi_1 u - u\|_2 + \|u\|_2 \leq Ch^{-1}\|u_h - \Pi_1 u\|_1 + C\|u\|_2 \\ &\leq Ch^{-1}(\|u - \Pi_1 u\|_1 + \|u - u_h\|_1) + C\|u\|_2 \\ &\leq Ch^{-1}\|u - u_h\|_1 + C\|u\|_2 \\ &\leq C(\|u\|_2 + \|f\|_2), \end{aligned}$$

where Theorem 4.5 is used in the last inequality. Therefore, we have

$$(4.30) \quad A(u - u_h, w_h) = \mathcal{O}(h^2)(\|f\|_2 + \|u\|_2)\|\theta\|_0.$$

For the second term (4.27), by continuity of the bilinear form and Lemma 4.4, we have

$$(4.31) \quad \begin{aligned} A(u - u_h, w - w_h) &\leq C\|u - u_h\|_1\|w - w_h\|_1 \leq Ch\|u - u_h\|_1\|w\|_2 \\ &\leq Ch\|u - u_h\|_1\|\theta\|_0 = \mathcal{O}(h^2)(\|f\|_2 + \|u\|_2)\|\theta\|_0. \end{aligned}$$

Therefore, by (4.27), (4.28) and (4.31), we have

$$(4.32) \quad \|\theta\|_0 = \mathcal{O}(h^2)(\|f\|_2 + \|u\|_2). \quad \square$$

REMARK 3. *Similar convergence results for the Q^1 method on general quasi-uniform quadrilateral meshes can be established via the same proof procedure in this section.*

384 **5. Extension to general quadrilateral meshes.** For a quadrilateral element
 385 e as in Fig. 2, let \mathbf{F}_e the mapping such that $\mathbf{F}_e(\hat{K}) = e$.

For $\varphi \in V_0^h$, by definition $\hat{\varphi} = \varphi|_e \circ \mathbf{F}_e \in Q^1(\hat{K})$. According to the chain rule, we have

$$\nabla \varphi \circ \mathbf{F}_e = DF_e^T \hat{\nabla} \hat{\varphi}$$

386 where $\varphi \circ \mathbf{F}_e = \hat{\varphi}$, $\nabla = \left(\frac{\partial}{\partial x_1}, \frac{\partial}{\partial x_2} \right)^T$, $\hat{\nabla} = \left(\frac{\partial}{\partial \hat{x}_1}, \frac{\partial}{\partial \hat{x}_2} \right)^T$.

387 Therefore, we have

$$388 \quad (5.1) \quad \int_e \mathbf{a} \nabla u_h \cdot \nabla v_h dx = \int_{\hat{K}} \left(DF_e^{-1} \hat{\mathbf{a}} DF_e^{T-1} \hat{\nabla} \hat{u}_h \right) \cdot \hat{\nabla} \hat{v}_h |J_e| d\hat{\mathbf{x}}$$

389 In the case of regular meshes with mesh size h , the matrix $DF_e^{-1} \hat{\mathbf{a}} DF_e^{*-1} = \frac{1}{h^2} \hat{\mathbf{a}}$
 390 and $J_e = h^2$.

391 Approximate (5.1) by the mixed quadrature (2.6) with parameter $\boldsymbol{\lambda} = (\lambda^1, \lambda^2)$,
 392 i.e.

$$393 \quad (5.2) \quad \int_e (\mathbf{a} \nabla u_h) \cdot \nabla v_h dx \approx \int_{\hat{K}} \left(\tilde{\mathbf{a}} \hat{\nabla} \hat{u}_h \right) \cdot \hat{\nabla} \hat{v}_h d_{\boldsymbol{\lambda}}^h \hat{\mathbf{x}}$$

394 where $\tilde{\mathbf{a}} = (|J_e| DF_e^{-1} \hat{\mathbf{a}} DF_e^{T-1}) \left(\frac{1}{2}, \frac{1}{2} \right)$.

As in Fig. 2, denote

$$\vec{\mathbf{c}}_0 = \mathbf{c}_{0,1} - \mathbf{c}_{0,0}, \quad \vec{\mathbf{c}}_1 = \mathbf{c}_{1,0} - \mathbf{c}_{0,0}, \quad \vec{\mathbf{c}}_2 = \mathbf{c}_{1,1} - \mathbf{c}_{1,0}, \quad \vec{\mathbf{c}}_3 = \mathbf{c}_{1,1} - \mathbf{c}_{0,1}$$

395 and

$$396 \quad \vec{\mathbf{c}}_i = (c_i^1, c_i^2)^T, \quad i = 0, 1, 2, 3, \quad DF_h = DF\left(\frac{1}{2}, \frac{1}{2}\right), \quad J_{e,h} = |J_e|\left(\frac{1}{2}, \frac{1}{2}\right), \quad \bar{\mathbf{a}}_e = \hat{\mathbf{a}}_e\left(\frac{1}{2}, \frac{1}{2}\right),$$

397 then we have

$$398 \quad DF_h = \frac{1}{2} \begin{pmatrix} c_1^1 + c_3^1 & c_0^1 + c_2^1 \\ c_1^2 + c_3^2 & c_0^2 + c_2^2 \end{pmatrix}, \quad DF_h^{-1} = \frac{1}{2 \det(DF_h)} \begin{pmatrix} c_0^2 + c_2^2 & -c_0^1 - c_2^1 \\ -c_1^2 - c_3^2 & c_1^1 + c_3^1 \end{pmatrix},$$

399
 400
 401

$$402 \quad (5.3) \quad \tilde{\mathbf{a}} = J_{e,h} DF_h^{-1} \bar{\mathbf{a}}_e DF_h^{T-1} = \begin{pmatrix} \tilde{a}_e^{11} & \tilde{a}_e^{12} \\ \tilde{a}_e^{12} & \tilde{a}_e^{22} \end{pmatrix}.$$

403 To make the stiffness matrix a M -matrix, by Theorem 3.2, the following is a
 404 sufficient condition:

$$405 \quad (5.4) \quad |\tilde{a}_e^{12}| \leq \min\{\tilde{a}_e^{11}, \tilde{a}_e^{22}\}.$$

406 While we have

$$407 \quad (5.5) \quad \begin{aligned} \tilde{a}^{11} &= \det(\bar{\mathbf{a}}_e) C \begin{pmatrix} c_0^2 + c_2^2 & -c_0^1 - c_2^1 \\ -c_1^2 - c_3^2 & c_1^1 + c_3^1 \end{pmatrix} \begin{pmatrix} \bar{a}^{11} & \bar{a}^{12} \\ \bar{a}^{12} & \bar{a}^{22} \end{pmatrix} \begin{pmatrix} c_0^2 + c_2^2 \\ -c_0^1 - c_2^1 \end{pmatrix} \\ &= \det(\bar{\mathbf{a}}_e) C \begin{pmatrix} c_0^1 + c_2^1 & c_0^2 + c_2^2 \\ c_1^1 + c_3^1 & c_1^2 + c_3^2 \end{pmatrix} \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} \bar{a}^{11} & \bar{a}^{12} \\ \bar{a}^{12} & \bar{a}^{22} \end{pmatrix} \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} \begin{pmatrix} c_0^1 + c_2^1 \\ c_0^2 + c_2^2 \end{pmatrix} \\ &= C \begin{pmatrix} c_0^1 + c_2^1 & c_0^2 + c_2^2 \\ c_1^1 + c_3^1 & c_1^2 + c_3^2 \end{pmatrix} \begin{pmatrix} \bar{a}^{22} & -\bar{a}^{12} \\ -\bar{a}^{12} & \bar{a}^{11} \end{pmatrix} \begin{pmatrix} c_0^1 + c_2^1 \\ c_0^2 + c_2^2 \end{pmatrix} \\ &= C (\vec{\mathbf{c}}_0 + \vec{\mathbf{c}}_2)^T \bar{\mathbf{a}}_e^{-1} (\vec{\mathbf{c}}_0 + \vec{\mathbf{c}}_2), \end{aligned}$$

408 and similarly

409 (5.6) $\tilde{a}^{12} = -C (\vec{\mathbf{c}}_0 + \vec{\mathbf{c}}_2)^T \bar{\mathbf{a}}_e^{-1} (\vec{\mathbf{c}}_1 + \vec{\mathbf{c}}_3), \quad \tilde{a}^{22} = C (\vec{\mathbf{c}}_1 + \vec{\mathbf{c}}_3)^T \bar{\mathbf{a}}_e^{-1} (\vec{\mathbf{c}}_1 + \vec{\mathbf{c}}_3),$

411 with $C = \frac{J_{e,h}}{4\det(DF_h)^2 \det(\bar{\mathbf{a}}_e)}$.

412 By $\vec{\mathbf{c}}_1 + \vec{\mathbf{c}}_2 - \vec{\mathbf{c}}_3 - \vec{\mathbf{c}}_0 = \vec{\mathbf{0}}$, (5.4) is equivalent to

413 (5.7)
$$\begin{aligned} (\vec{\mathbf{c}}_0 + \vec{\mathbf{c}}_2)^T \bar{\mathbf{a}}_e^{-1} (\vec{\mathbf{c}}_0 + \vec{\mathbf{c}}_3) &> 0, & (\vec{\mathbf{c}}_0 + \vec{\mathbf{c}}_2)^T \bar{\mathbf{a}}_e^{-1} (\vec{\mathbf{c}}_0 - \vec{\mathbf{c}}_1) &> 0, \\ (\vec{\mathbf{c}}_1 + \vec{\mathbf{c}}_3)^T \bar{\mathbf{a}}_e^{-1} (\vec{\mathbf{c}}_0 + \vec{\mathbf{c}}_3) &> 0, & (\vec{\mathbf{c}}_1 + \vec{\mathbf{c}}_3)^T \bar{\mathbf{a}}_e^{-1} (\vec{\mathbf{c}}_1 - \vec{\mathbf{c}}_0) &> 0. \end{aligned}$$

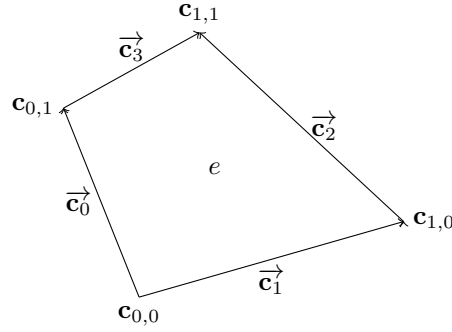


FIG. 2. A quadrilateral element e .

414 **THEOREM 5.1.** *If the quadrilateral mesh fulfill the condition (5.4) with $\tilde{\mathbf{a}}$ defined*
 415 *in (5.3) or the mesh condition (5.7), then the stiffness matrix of the linear Q^1 finite*
 416 *element scheme (3.2) for solving BVP (1.1) is an M-matrix.*

417 **REMARK 4.** *If the diffusion operator degenerate to Laplacian, i.e. $\mathbf{a} = \alpha(\mathbf{x})I$. A*
 418 *sufficient condition for (5.7) is that both diagonals of the quadrilateral element bisect*
 419 *each angle, resulting in two non-obtuse angles for each vertex.*

420 **REMARK 5.** *By adopting some anisotropic mesh adaptation strategy where an*
 421 *anisotropic mesh is generated as an M-uniform mesh or a uniform mesh in the metric*
 422 *specified by the diffusion matrix \mathbf{a} . The method (3.2) for any anisotropic problem*
 423 *possibly can be monotone on that anisotropic mesh.*

424 If we consider rectangular meshes, for simplicity we assume

425
$$\mathbf{c}_{0,0} = (0, 0), \quad \mathbf{c}_{1,0} = (h_1, 0), \quad \mathbf{c}_{0,1} = (0, h_2), \quad \mathbf{c}_{1,1} = (h_1, h_2).$$

426 Then we have

427
$$\tilde{\mathbf{a}} = \begin{pmatrix} \frac{h_2}{h_1} \bar{a}^{11} & \bar{a}^{12} \\ \bar{a}^{12} & \frac{h_1}{h_2} \bar{a}^{22} \end{pmatrix}$$

428 and (5.4) becomes

429 (5.8)
$$|\bar{a}_e^{12}| \leq \min\left\{\frac{h_2}{h_1} \bar{a}_e^{11}, \frac{h_1}{h_2} \bar{a}_e^{22}\right\}.$$

430 Recall that $\sqrt{\bar{a}_e^{11}\bar{a}_e^{22}} \geq |\bar{a}_e^{12}|$, taking $\frac{h_1}{h_2} = \sqrt{\frac{\bar{a}_e^{11}}{\bar{a}_e^{22}}}$ will guarantee (5.8). Therefore, if the
 431 rectangular mesh is deployed with aspect ratio $\sqrt{\frac{\bar{a}_e^{11}}{\bar{a}_e^{22}}}$, then the stiffness matrix of the
 432 Q^1 method (3.2) is a M -matrix.

433 If the elliptic coefficient \mathbf{a} is constant on the whole domain Ω , when the rect-
 434 angular mesh are fine enough, there must exist rectangular mesh with aspect ratio
 435 approximately $\sqrt{\frac{\bar{a}_e^{11}}{\bar{a}_e^{22}}}$ such that the stiffness matrix of scheme (3.2) solve the BVP (1.1)
 436 is an M -matrix.

437 **6. Numerical experiment.** In this section, we show an accuracy test verifying
 438 the proved order of accuracy of the scheme (3.2) on uniform meshes. We consider the
 439 following two dimensional elliptic equation:

$$440 \quad (6.1) \quad -\nabla \cdot (\mathbf{a}\nabla u) + cu = f \quad \text{on } [0, \pi]^2$$

where $\mathbf{a} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$, $a_{11} = a_{12} = a_{21} = 1 + 10x_2^2 + x_1 \cos x_2 + x_2$, $a_{22} =$
 $2 + 10x_2^2 + x_1 \cos x_2 + x_2$, with an exact solution

$$u(x_1, x_2) = -\sin x_1^2 \sin x_2 \cos x_2.$$

441 The errors at grid points are listed in Table 1. We observe the desired second
 442 order accuracy in the discrete 2-norm and infinity norm for the function values.

TABLE 1

A 2D elliptic equation with Dirichlet boundary conditions. The first column is the number of
 elements in a finite element mesh. The second column is the number of degree of freedoms.

FEM Mesh	DoF	l^2 error	order	l^∞ error	order
4×4	3^2	4.41E-1	-	3.48E-1	-
8×8	7^2	7.20E-2	2.61	5.93E-2	2.55
16×16	15^2	1.65E-2	2.13	1.39E-2	2.09
32×32	31^2	4.03E-3	2.03	3.45E-3	2.02
64×64	63^2	1.00E-3	2.01	8.61E-4	2.00

443 **7. Conclusion.** We constructed a linear monotone Q^1 finite element method
 444 for anisotropic diffusion problem (1.1). On uniform meshes, when the diffusion matrix
 445 is diagonally dominant, the M -matrix property is guaranteed thus monotonicity is
 446 achieved. When this Q^1 finite element method is deployed on a general quadrilateral
 447 mesh, we get a local mesh constraint.

- 449 [1] J. BRAMBLE, B. HUBBARD, AND V. THOMÉE, *Convergence estimates for essentially positive*
 450 *type discrete dirichlet problems*, Mathematics of Computation, 23 (1969), pp. 695–709.
 451 [2] C. CANCÈS, M. CATHALA, AND C. LE POTIER, *Monotone corrections for generic cell-centered*
 452 *finite volume approximations of anisotropic diffusion equations*, Numerische Mathematik,
 453 125 (2013), pp. 387–417.
 454 [3] I. CHRISTIE AND C. HALL, *The maximum principle for bilinear elements*, Internat. J. Numer.
 455 Methods Engrg., 20 (1984), pp. 549–553.
 456 [4] P. G. CIARLET, *The finite element method for elliptic problems*, Classics in applied mathemat-
 457 ics, 40 (2002), pp. 1–511.

- 458 [5] L. J. CROSS AND X. ZHANG, *Monotonicity of Q^3 spectral element method for discrete Laplacian*,
459 2023, <https://arxiv.org/abs/2010.07282>.
- 460 [6] L. J. CROSS AND X. ZHANG, *On the monotonicity of Q^2 spectral element method for Laplacian*
461 *on quasi-uniform rectangular meshes*, to appear in *Communications in Computational*
462 *Physics*, (2023).
- 463 [7] L. C. EVANS, *Partial Differential Equations*, vol. 019 of Graduate Studies in Mathematics,
464 American Mathematical Society, 2 ed., 2010.
- 465 [8] D. GREENSPAN AND P. JAIN, *On non negative difference analogues of elliptic differential equa-*
466 *tions*, *Journal of the Franklin Institute*, 279 (1965), pp. 360–365.
- 467 [9] J. HU AND X. ZHANG, *Positivity-preserving and energy-dissipative finite difference schemes*
468 *for the Fokker–Planck and Keller–Segel equations*, *IMA Journal of Numerical Analysis*, 43
469 (2022), pp. 1450–1484.
- 470 [10] W. HUANG, *Discrete maximum principle and a delaunay-type mesh condition for linear fi-*
471 *nite element approximations of two-dimensional anisotropic diffusion problems*, *Numerical*
472 *Mathematics: Theory, Methods and Applications*, 4 (2011), pp. 319–334.
- 473 [11] D. KUZMIN, M. J. SHASHKOV, AND D. SVYATSKIY, *A constrained finite element method satisfy-*
474 *ing the discrete maximum principle for anisotropic diffusion problems*, *J. Comput. Phys.*,
475 228 (2009), pp. 3448–3463.
- 476 [12] C. LE POTIER, *A nonlinear finite volume scheme satisfying maximum and minimum principles*
477 *for diffusion operators*, *International Journal on Finite Volumes*, (2009), pp. 1–20.
- 478 [13] H. LI, S. XIE, AND X. ZHANG, *A high order accurate bound-preserving compact finite difference*
479 *scheme for scalar convection diffusion equations*, *SIAM Journal on Numerical Analysis*,
480 56 (2018), pp. 3308–3345.
- 481 [14] H. LI AND X. ZHANG, *On the monotonicity and discrete maximum principle of the finite differ-*
482 *ence implementation of $C^0 C^0$ - $Q^2 Q^2$ finite element method*, *Numerische Mathematik*,
483 145 (2020), pp. 437–472.
- 484 [15] H. LI AND X. ZHANG, *Superconvergence of high order finite difference schemes based on varia-*
485 *tional formulation for elliptic equations*, *Journal of Scientific Computing*, 82 (2020), pp. 1–
486 39.
- 487 [16] H. LI AND X. ZHANG, *A high order accurate bound-preserving compact finite difference scheme*
488 *for two-dimensional incompressible flow*, *Communications on Applied Mathematics and*
489 *Computation*, (2023), pp. 1–29.
- 490 [17] X. LI AND W. HUANG, *An anisotropic mesh adaptation method for the finite element solution*
491 *of heterogeneous anisotropic diffusion problems*, *Journal of Computational Physics*, 229
492 (2010), pp. 8072–8094.
- 493 [18] X. LI, D. SVYATSKIY, AND M. SHASHKOV, *Mesh adaptation and discrete maximum principle*
494 *for 2d anisotropic diffusion problems*, tech. report, Technical Report LA-UR 10-01227, Los
495 Alamos National Laboratory, Los Alamos, NM, 2007.
- 496 [19] K. LIPNIKOV, M. SHASHKOV, D. SVYATSKIY, AND Y. VASSILEVSKI, *Monotone finite volume*
497 *schemes for diffusion equations on unstructured triangular and shape-regular polygonal*
498 *meshes*, *Journal of Computational Physics*, 227 (2007), pp. 492–512.
- 499 [20] R. LISKA AND M. SHASHKOV, *Enforcing the discrete maximum principle for linear finite element*
500 *solutions of second-order elliptic problems*, *Commun. Comput. Phys.*, 3 (2008), pp. 852–
501 877.
- 502 [21] C. LIU, Y. GAO, AND X. ZHANG, *Structure preserving schemes for fokker-planck equations of*
503 *irreversible processes*, to appear in *Journal of Scientific Computing*, (2023).
- 504 [22] C. LIU AND X. ZHANG, *A positivity-preserving implicit-explicit scheme with high order poly-*
505 *nomial basis for compressible Navier–Stokes equations*, *Journal of Computational Physics*,
506 493 (2023), p. 112496.
- 507 [23] J. LORENZ, *Zur inversmonotonie diskreter probleme*, *Numer. Math.*, 27 (1977), pp. 227–238.
- 508 [24] C. LU, W. HUANG, AND J. QIU, *Maximum principle in linear finite element approximations of*
509 *anisotropic diffusion–convection–reaction problems*, *Numerische Mathematik*, 127 (2014),
510 pp. 515–537.
- 511 [25] M. J. MLACNIK AND L. J. DURLOFSKY, *Unstructured grid optimization for improved monotonic-*
512 *ity of discrete solutions of elliptic equations with highly anisotropic coefficients*, *Journal of*
513 *Computational Physics*, 216 (2006), pp. 337–361.
- 514 [26] T. S. MOTZKIN AND W. WASOW, *On the approximation of linear elliptic differential equations*
515 *by difference equations with positive coefficients*, *Journal of Mathematics and Physics*, 31
516 (1952), pp. 253–259.
- 517 [27] C. NGO AND W. HUANG, *Monotone finite difference schemes for anisotropic diffusion prob-*
518 *lems via nonnegative directional splittings*, *Communications in Computational Physics*, 19
519 (2016), pp. 473–495.

- 520 [28] J. M. NORDBOTTEN, I. AAVATSMARK, AND G. EIGESTAD, *Monotonicity of control volume meth-*
521 *ods*, *Numerische Mathematik*, 106 (2007), pp. 255–288.
- 522 [29] R. J. PLEMMONS, *M-matrix characterizations. I-nonsingular M-matrices*, *Numer. Anal. Appl.*,
523 18 (1977), pp. 175–188.
- 524 [30] P. SHARMA AND G. W. HAMMETT, *Preserving monotonicity in anisotropic diffusion*, *Journal*
525 *of Computational Physics*, 227 (2007), pp. 123–142.
- 526 [31] J. SHEN AND X. ZHANG, *Discrete maximum principle of a high order finite difference scheme for*
527 *a generalized Allen–Cahn equation*, *Communications in Mathematical Sciences*, 20 (2022),
528 pp. 1409–1436.
- 529 [32] J. WEICKERT ET AL., *Anisotropic diffusion in image processing*, vol. 1, Teubner Stuttgart, 1998.
- 530 [33] J. XU AND L. ZIKATANOV, *A monotone finite element scheme for convection-diffusion equa-*
531 *tions*, *Math. Comp.*, 68 (1999), pp. 1429–1446.
- 532 [34] G. YUAN AND Z. SHENG, *Monotone finite volume schemes for diffusion equations on polygonal*
533 *meshes*, *Journal of computational physics*, 227 (2008), pp. 6288–6312.