12

A MONOTONE Q¹ FINITE ELEMENT METHOD FOR ANISOTROPIC ELLIPTIC EQUATIONS

3

HAO LI * and XIANGXIONG ZHANG †

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

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)
$$\mathcal{L}u \equiv -\nabla \cdot (\mathbf{a}\nabla u) + cu = f \quad \text{on} \quad \Omega,$$
$$u = g \quad \text{on} \quad \partial\Omega,$$

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 $\overline{\Omega}$ 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 $u \in C^2(\Omega) \cap C(\overline{\Omega})$, a maximum principle holds [8]:

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

16 In particular,

17 (1.2)
$$\mathcal{L}u = 0 \text{ in } \Omega \Longrightarrow |u(\mathbf{x})| \le \max_{\partial \Omega} |u|, \quad \forall (\mathbf{x}) \in \Omega.$$

The anisotropic diffusion problem (1.1) arises from various areas of science and 18 engineering, including plasma physics, Lagrangian hydrodynamics, and image pro-19 cessing. To avoid spurious oscillations or non-physical numerical solution, it is de-20sired to have numerical schemes to satisfy (1.2) in the discrete sense. We are in-21terested in a linear approximation to \mathcal{L} which can be represented as a matrix L_h . 22 The matrix L_h is called monotone if its inverse only has nonnegative entries, i.e., 23 $L_h^{-1} \ge 0$. Monotonicity of the scheme is a sufficient condition for the discrete max-24 imum principle and has various applications espeically for parabolic problems, see 25[1, 34, 15, 10, 32, 22, 7, 6, 23, 22, 14, 17].26

^{*}Department of Applied Mathematics, The Hong Kong Polytechnic University, Hong Kong. Email: hao94.li@polyu.edu.hk.

[†]Department of Mathematics, Purdue University, 150 N. University Street, West Lafayette, IN 47907-2067, USA. Email: zhan1966@purdue.edu.

H. LI AND X. ZHANG

1.2. Monotone schemes for anisotropic diffusion equations. Monotone (or positive-type in some literature) numerical methods for problem (1.1) have received considerable attention, e.g., see [12, 18, 19, 20, 21, 26, 35, 31, 13, 2, 28]. The major efforts of studying linear monotone schemes take advantage of *M*-matrix (see [30] for the definition), either by showing the coefficient matrix is *M*-matrix directly or the coefficient matrix can be factorized into a product of *M*-matrices. In the following, we call a numerical scheme satisfying *M*-matrix property if the corresponding coefficient matrix is an *M*-matrix.

By factorizing the stiffness matrix into a product of *M*-matrices, the monotonocity can still be ensured. For a nine-point scheme on a two-dimensional quadrilateral grid, the matrix condition for monotonicity with specific splitting strategy in [29] aligns with the Lorenz's condition presented in [24, 15]. The difference is that in [24, 15], only the existence of the factorization was proved while in [29] the exact matrix factorization was found explicitly.

In [27], it is proved that a monotone finite difference scheme exists for any lin-41 ear second-order elliptic problem on fine enough uniform mesh but a finite difference 42 method with fixed stencil for all the problems satisfying the M-matrix property does 43 not exist. With nonnegative directional splittings, [33, 9, 28] propose to construct 44 finite difference schemes for elliptic operators in the non-divergence form and diver-45 gence form. Particularly in [28], it is shown that a monotone scheme satisfying the 46 *M*-matrix property can be constructed for continuous diffusion matrix for sufficiently 47 fine mesh and sufficiently large finite difference stencil. 48

In [18], for the P^1 finite elements in two and three dimensions, the author gen-49 eralized the well known non-obtuse angle condition for anisotropic diffusion problem 50in the sense to have the dihedral angles of all mesh elements, measured in a metric 51depending on $\mathbf{a}(\mathbf{x})$, be non-obtuse. It reduces to the non-obtuse angle condition for isotropic diffusion matrices when $\mathbf{a}(\mathbf{x}) = \alpha(\mathbf{x})\mathbb{I}$, where \mathbb{I} is the identity matrix. The 53 formulation was also utilized in [18] for the construction of the so called *M*-uniform 5455meshes on which the numerical scheme is monotone. The approach to show monotonicity in [18] is to write the global matrix as the sum of local contributions. In 56 [11], the Delaunay condition is extended to anisotropic diffusion problems through a 57 refined analysis studying the whole stiffness matrix for the two-dimensional situation. 58The analysis of [18] was extended to the anisotropic diffusion-convection-reaction problems in [25]. 60

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

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

To analyze the monotonicity of the stiffness matrix, we approximate integrals with a specific quadrature rule, particularly, the linear combination of the trapezoid rule and midpoint rule. Then we demonstrate that the continuous Q^1 finite element

MONOTONE Q^1 FEM

⁷⁶ method with the specific quadrature rule, when applied to the anisotropic diffusion

77 problem on a uniform mesh, ensures monotonicity for the problem with a diagonally

78 dominant diffusion coefficient matrix. The method is linear and second-order accurate.

The convergence of the function values for this method is rigorously proven. The coefficient constraints become mesh constraints when this Q^1 finite element method is applied on general quadrilateral meshes.

This paper is organized as follows. In Section 2, we introduce the notations and review some standard quadrature error estimates. In Section 3, we derive the Q^1 scheme for anisotropic diffusion equation with Dirichlet boundary condition and the coefficient constraints for the stiffness matrix to be an *M*-matrix. In Section 4, the convergence rate of function values is proved. In Section 5, we discuss the extension to general quadrilateral meshes. Numerical results are given in Section 6.

88 2. Preliminaries.

89 2.1. Notation and tools. We introduce some notation and useful tools as fol-90 lows.

- For the problem dimension d, though we only consider the case d = 2, sometimes we keep the general notation d to illustrate how the results are influenced by the dimension.
- For the Q¹ finite element space, i.e., tensor product of linear polynomials, the local space is defined on a reference element K, e.g., K = [0,1]². Then, the finite element space on a physical mesh element e is given by the reference map from K to e. The reference element K is as Figure 1.



FIG. 1. The reference element.

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

(2.1)

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

101 102

100

97

• We will use \wedge for a function to emphasize the function is defined on or transformed to the reference element \hat{K} from a physical mesh element.

• For a quadrilateral element e, we assume $\mathbf{F}_e = (F_{e1}, F_{e2})^T$ 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

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

103 104 • $Q^1(\hat{K}) = \left\{ p(\mathbf{x}) = \sum_{i=0}^1 \sum_{j=1}^1 p_{ij} \hat{\phi}_{i,j}(\hat{\mathbf{x}}), \, \hat{\mathbf{x}} \in \hat{K} \right\}$ is the set of Q^1 polynomials on the reference element \hat{K} .

H. LI AND X. ZHANG

• Inverse estimates for polynomials: there exists a constant $C_I > 0$, independent of h and e, such that for

$$||v_h||_{k+1} \le C_I h^{-1} ||v_h||_k, \quad \forall v_h \in V^h, k \ge 0$$

• Elliptic regularity holds for the problem (3.1):

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

138 • Let Ω_h is a finite element mesh for Ω . For each element $e \in \Omega_h$, we denote 139 $\bar{\mathbf{a}}_e = (\bar{a}_e^{ij})$ as an approximation to the average of \mathbf{a} on element e, i.e. $\bar{a}_e^{ij} = \frac{1}{meas(e)} \int_e a^{ij} d\mathbf{x}$. Then we define piece-wise constant function $\bar{\mathbf{a}}$ on Ω as

141
$$\bar{\mathbf{a}}(\mathbf{x}) = \bar{\mathbf{a}}_e, \text{ for } \mathbf{x} \in e.$$

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

143 (2.2)
$$\int_{\hat{K}} \left(\hat{\Pi}_1 \hat{u} \right) \hat{w} d\hat{\mathbf{x}} = \int_{\hat{K}} \hat{u} \hat{w} d\hat{\mathbf{x}}, \quad \forall \hat{w} \in Q^1(\hat{K}).$$

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:

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

From which it follows that:

146

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

144 This establishes that $\hat{\Pi}_1$ acts as a continuous linear mapping from $L^2(\hat{K})$ to 145 $H^1(\hat{K})$. Similarly, by extending this argument, we can also demonstrate that

- $H^1(\hat{K})$. 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 , $j = 1, \ldots, N_h$ and all the the vertices of Ω_h on $\partial\Omega$ by \mathbf{x}_j , $j = N_h + 1, \ldots, N_h + N_h^\partial$. The corresponding Lagrange basis functions in V_h are denoted by φ_i , $i = 1, \ldots, N_h + N_h^\partial$, which are continuous in Ω , linear polynomials in each element e and

$$\varphi_i(\mathbf{x}_i) = \delta_{ij}, \quad i, j = 1, \dots, N_h + N_h^{\partial}.$$

147 **2.2. Mixed quadrature.** To analyze and impose the monotonicity of the stiff-148 ness matrix, we will use numerical quadrature rules to approximate integrals. As we 149 will see, the choice of quadrature rules can significantly affect the monotonicity of the 150 numerical schemes.

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)
$$\int_{0}^{1} f(\hat{x}) d\hat{x} \simeq \lambda \frac{f(0) + f(1)}{2} + (1 - \lambda) f\left(\frac{1}{2}\right) \\ = \hat{\omega}_{1} f(\hat{\xi}_{1}) + \hat{\omega}_{2} f(\hat{\xi}_{2}) + \hat{\omega}_{3} f(\hat{\xi}_{1}),$$

156 where λ is a parameter to be determined and

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

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,

$$\int_{\hat{K}} f(\hat{\mathbf{x}}) d\hat{\mathbf{x}} = \int_{0}^{1} \int_{0}^{1} f(\hat{x}_{1}, \hat{x}_{2}) d\hat{x}_{1} d\hat{x}_{2} = \int_{0}^{1} \left(\int_{0}^{1} f(\hat{x}_{1}, \hat{x}_{2}) d\hat{x}_{2} \right) d\hat{x}_{1}$$

$$\simeq \int_{0}^{1} \left(\sum_{q=1}^{3} \hat{\omega}_{q}^{2} f\left(\hat{x}_{1}, \hat{\xi}_{q}\right) \right) d\hat{x}_{1} \simeq \sum_{p=1}^{3} \hat{\omega}_{p}^{1} \left(\sum_{q=1}^{3} \hat{\omega}_{q}^{2} f\left(\hat{\xi}_{p}, \hat{\xi}_{q}\right) \right) = \sum_{p=1}^{3} \sum_{q=1}^{3} \hat{\omega}_{p}^{1} \hat{\omega}_{q}^{2} f\left(\hat{\xi}_{p}, \hat{\xi}_{q}\right),$$

164 where ω_i^j are just ω_i but 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 (2.5)
$$\int_{\hat{K}} \hat{f}(\hat{\mathbf{x}}) d^h_{\boldsymbol{\lambda}} \hat{\mathbf{x}} := \sum_{p=1}^3 \sum_{q=1}^3 \hat{\omega}_p^1 \hat{\omega}_q^2 f\left(\hat{\xi}_p, \hat{\xi}_q\right).$$

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

171 (2.6)
$$\int_{e} f(\mathbf{x}) d^{h}_{\boldsymbol{\lambda}_{e}} \mathbf{x} := \int_{\hat{K}} f \circ \mathbf{F}_{e}(\hat{\mathbf{x}}) d^{h}_{\boldsymbol{\lambda}_{e}} \hat{\mathbf{x}}.$$

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

173 (2.7)
$$\int_{\Omega} f d^{h}_{\lambda_{\Omega}} \mathbf{x} := \sum_{e \in \Omega_{h}} \int_{e} f d^{h}_{\lambda_{e}} \mathbf{x}$$

where $\lambda_{\Omega} = (\lambda_e)_{e \in \Omega_h}$ can be viewed as a vector-valued piece-wise constant function, with values λ_e which differ across different elements.

As a particular instance, $\int_{\Omega} f d_1^h \mathbf{x}$ denotes the case $\lambda_e = (1,0)$ for all $e \in \Omega_h$, i.e. 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 [5]:

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

182 (2.8)
$$\|\hat{u} - \hat{\Pi}\hat{u}\|_{k+1,\hat{K}} \le C[\hat{u}]_{k+1,\hat{K}}, \quad \forall \hat{u} \in H^{k+1}(\hat{K}).$$

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

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

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

(2.9)
$$\int_{e} a d\mathbf{x} - \int_{e} a d^{h} \mathbf{x} = \mathcal{O}\left(h^{2+\frac{d}{2}}\right) [a]_{2,e} = \mathcal{O}\left(h^{2+d}\right) [a]_{2,\infty,e}$$

190 (2.10)
$$\int_{e} a d\mathbf{x} - \int_{e} \bar{a}_{e} d\mathbf{x} = \mathcal{O}\left(h^{2+\frac{d}{2}}\right) [a]_{2,e} = \mathcal{O}\left(h^{2+d}\right) [a]_{2,\infty,e}$$

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

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

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

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

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

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

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

193

195

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

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

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

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

$$\int_{e} u_{x_{i}}(v_{h})_{x_{j}} d\mathbf{x} - \int u_{x_{i}}(v_{h})_{x_{j}} d^{h}_{\boldsymbol{\lambda}_{e}} \mathbf{x} = \mathcal{O}\left(h^{2}\right) \|u\|_{3,e} \|v_{h}\|_{2,e}.$$

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

$$\begin{split} E(u_{x_{i}}(v_{h})_{x_{j}}) &= h^{d-2} \hat{E}(\hat{u}_{\hat{x}_{i}}(\hat{v}_{h})_{\hat{x}_{j}}) \leq Ch^{d-2} [\hat{u}_{\hat{x}_{i}}(\hat{v}_{h})_{\hat{x}_{j}}]_{2,2,\hat{K}} \\ &\leq Ch^{d-2} \left(|\hat{u}_{\hat{x}_{i}}|_{2,2,\hat{K}} |(\hat{v}_{h})_{\hat{x}_{j}}|_{0,\infty,\hat{K}} + |\hat{u}_{\hat{x}_{i}}|_{1,2,\hat{K}} |(\hat{v}_{h})_{\hat{x}_{j}}|_{1,\infty,\hat{K}} + |\hat{u}_{\hat{x}_{i}}|_{0,2,\hat{K}} |(\hat{v}_{h})_{\hat{x}_{j}}|_{2,\infty,\hat{K}} \right) \\ &\leq Ch^{d-2} \left(|\hat{u}_{\hat{x}_{i}}|_{2,2,\hat{K}} |(\hat{v}_{h})_{\hat{x}_{j}}|_{0,2,\hat{K}} + |\hat{u}_{\hat{x}_{i}}|_{1,2,\hat{K}} |(\hat{v}_{h})_{\hat{x}_{j}}|_{1,2,\hat{K}} + |\hat{u}_{\hat{x}_{i}}|_{0,2,\hat{K}} |(\hat{v}_{h})_{\hat{x}_{j}}|_{2,2,\hat{K}} \right) \\ &\leq Ch^{d-2} \left(|\hat{u}|_{3,2,\hat{K}} |\hat{v}_{h}|_{1,2,\hat{K}} + |\hat{u}|_{2,2,\hat{K}} |\hat{v}_{h}|_{2,2,\hat{K}} \right). \end{split}$$

where the second last inequality is implied by the equivalence of norms over $Q^1(\hat{K})$

and in the last inequality we use the fact that the third derivative of Q^1 polynomial vanish. 199 Therefore,

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

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

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

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:

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

3. The Q^1 finite element method and its monotonicity. In this section, we first derive the Q^1 finite element scheme then pursue its monotonicity.

3.1. Derivation of the scheme. For problem (1.1), assuming there is a function $\bar{g} \in H^1(\Omega)$ as an extension of g so that $\bar{g}|_{\partial\Omega} = g$, the variational form of (1.1) is to find $\tilde{u} = u - \bar{g} \in H^1_0(\Omega)$ satisfying

208 (3.1)
$$\mathcal{A}(u,v) = (f,v) - \mathcal{A}(\bar{g},v), \quad \forall v \in H_0^1(\Omega).$$

209 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 $\frac{1}{meas(e)} \int_e \mathbf{a} d\mathbf{x}$ or its middle point value on each element e. The approximation is denoted by $\bar{\mathbf{a}}_e$. Then we obtain the modified bilinear form

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

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

By approximating integrals in $\bar{\mathcal{A}}(u, v)$ with quadrature specified in (2.7), along with designated quadrature parameter λ_{Ω} , we derive the following numerical scheme: find $u_h \in V_0^h$ satisfying

215 (3.2)
$$\mathcal{A}_h(u_h, v_h) = \langle f, v_h \rangle_h - A_h(g_I, v_h), \quad \forall v_h \in V_0^h,$$

216 where the approximated bilinear form is defined as

217 (3.3)
$$\mathcal{A}_h(u_h, v_h) := \int_{\Omega} \bar{\mathbf{a}} \nabla u_h \cdot \nabla v_h d^h_{\lambda} \mathbf{x} + \int_{\Omega} c u_h v_h d^h_1 \mathbf{x}.$$

218 The right hand side is defined as

219 (3.4)
$$\langle f, v_h \rangle_h := \int_{\Omega} f v_h d_1^h \mathbf{x},$$

This manuscript is for review purposes only.

8

20

and $g_I \in V^h$ is the piece-wise Q^1 Lagrangian interpolation polynomial of the following function:

$$g(x,y) = \begin{cases} 0, & \text{if } (x,y) \in (0,1)^2\\ g(x,y), & \text{if } (x,y) \in \partial\Omega. \end{cases}$$

220 Then $\bar{u}_h = u_h + g_I$ is the numerical solution for the problem (1.1).

Obviously the quadrature parameters $\boldsymbol{\lambda} = (\lambda^1, \lambda^2)$ on each element are to be determined for the quadrature (2.6). It is not obvious that the numerical solution \bar{u}_h is an accurate approximation of the exact solution u as $\bar{\mathbf{a}}$ varies depending on the mesh.

Let us denote **f** the vector consisting of $f_i = f(\mathbf{x}_i)$ for $i = 1, ..., N_h$ and \mathbf{f} an abstract vector consisting of f_i for $i = 1, ..., N_h$ and the boundary condition $g_i = g(\mathbf{x}_i)$ at the boundary grid points $i = N_h + 1, ..., N_h + N_h^{\partial}$. Besides, we denote $\mathbf{\bar{u}} = (u_1, ..., u_{N_h + N_h^{\partial}})$ the vector such that

$$\bar{u}_h = \sum_{i=1}^{N_h + N_h^\partial} u_i \varphi_i$$

Then scheme (3.2) can be written as a finite difference scheme [16], with the matrix vector representation $\bar{A}\bar{\mathbf{u}} = M\mathbf{f}$ where $\bar{A} = (a_{ij})_{N_h \times (N_h + N_h^\partial)}$, $a_{ij} = \mathcal{A}_h(\varphi_j, \varphi_i)$, i =1,..., N_h , $j = 1, ..., N_h + N_h^\partial$, and M is the lumped mass matrix. For convenience, after inverting the lumped mass matrix M, with the boundary conditions, the whole scheme can be represented in a matrix vector form

230 (3.5)
$$\bar{L}_h \bar{\mathbf{u}} = \bar{\mathbf{f}},$$

231 where

232
$$(\bar{L}_h \bar{\mathbf{u}})_i := (M^{-1} \bar{A} \bar{\mathbf{u}})_i = f_i, \quad i = 1, \dots, N_h,$$

$$(\bar{L}_h \bar{\mathbf{u}})_i := u_i = g_i, \quad i = N_h + 1, \dots, N_h + N_h^{\partial}.$$

3.2. Discrete maximum principle. In this subsection, we review how the monotonicity implies the discrete maximum principle. For the matrix form (3.5) of the scheme (3.2), with

$$\mathbf{u} = (u_1, \dots, u_{N_h})^T, \quad \mathbf{u}^\partial = \left(u_{N_h+1}, \dots, u_{N_h+N_h^\partial}\right)^T, \quad \bar{\mathbf{u}} = \left(u_1, \dots, u_{N_h+N_h^\partial}\right)^T,$$

we have the finite difference operator \mathcal{L}_h defined by \bar{L}_h

$$\mathcal{L}_{h}(\bar{\mathbf{u}}) := \bar{L}_{h}\bar{\mathbf{u}} = \bar{\mathbf{f}}, \quad \bar{L}_{h} = \begin{pmatrix} L_{h} & B^{\partial} \\ 0 & I \end{pmatrix}, \ \bar{\mathbf{u}} = \begin{pmatrix} \mathbf{u} \\ \mathbf{u}^{\partial} \end{pmatrix}, \ \bar{\mathbf{f}} = \begin{pmatrix} \mathbf{f} \\ \mathbf{g} \end{pmatrix}$$

235 The discrete maximum principle is

236 (3.6)
$$\bar{L}_h(\bar{\mathbf{u}})_i \le 0, \ 1 \le i \le N_h \Longrightarrow \max_i u_i \le \max\left\{0, \max_{N_h+1 \le i \le N_h+N_h^\partial} u_i\right\}.$$

The following result was proven in [4]:

THEOREM 3.1. A finite difference operator \mathcal{L}_h satisfies the discrete maximum principle (3.6) if $\bar{L}_h^{-1} \geq 0$ and all row sums of \bar{L}_h are non-negative. 3.3. Monotonicity of the Q^1 finite element. To have the monotonicity, we are interested in conditions for \bar{L}_h being an *M*-matrix. Recall a sufficient condition for *M*-matrix, see condition C_{10} in [30]:

LEMMA 3.2. For a real irreducible square matrix A with positive diagonal entries and non-positive off-diagonal entries, A is a nonsingular M-matrix if all the row sums of A are non-negative and at least one row sum is positive.

Then we have the following result on the uniform rectangular mesh. The stiffness matrix of (3.2) is denoted as $A = (a_{ij}) = (\mathcal{A}_h(\varphi_j, \varphi_i)), i, j = 1, \dots, N_h$.

THEOREM 3.3. 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 given by (3.2) for the elliptic equation (1.1) on uniform rectangular mesh, the stiffness matrix and \bar{L}_h are *M*-matrices and the finite difference operator defined by \bar{L}_h satisfies discrete maximum principle, provided the quadrature parameters for each element e are chosen as:

253 (3.7)
$$\lambda_e^1, \, \lambda_e^2 \in \left(\frac{|\bar{a}_e^{11} - \bar{a}_e^{22}|}{\bar{a}_e^{11} + \bar{a}_e^{22}}, 1 - \frac{2|\bar{a}_e^{12}|}{\bar{a}_e^{11} + \bar{a}_e^{22}}\right].$$

254 When $|\bar{a}_e^{12}| = \min\{\bar{a}_e^{11}, \bar{a}_e^{22}\}$, (3.7) means we take λ_e^1, λ_e^2 to be the upper bound of the 255 interval, i.e. $1 - \frac{2|\bar{a}_e^{12}|}{\bar{a}_e^{11} + \bar{a}_e^{22}}$.

256 *Proof.* First, we consider the following quadrature approximation results on the 257 reference element \hat{K} . With quadrature (2.5) and quadrature parameter $\lambda_e = (\lambda_e^1, \lambda_e^2)$, 258 we have

259
$$\langle \bar{\mathbf{a}} \nabla \phi_{0,0}, \nabla \phi_{0,1} \rangle_h = \langle \bar{\mathbf{a}} \nabla \phi_{1,1}, \nabla \phi_{1,0} \rangle_h = -\frac{1}{4} (\lambda_e^2 \bar{a}_e^{11} + \lambda_e^1 \bar{a}_e^{22}) + \frac{1}{4} (\bar{a}_e^{11} - \bar{a}_e^{22}),$$

260
$$\langle \bar{\mathbf{a}} \nabla \phi_{0,0}, \nabla \phi_{1,0} \rangle_h = \langle \bar{\mathbf{a}} \nabla \phi_{0,1}, \nabla \phi_{1,1} \rangle_h = -\frac{1}{4} (\lambda_e^2 \bar{a}_e^{11} + \lambda_e^1 \bar{a}_e^{22}) + \frac{1}{4} (\bar{a}_e^{22} - \bar{a}_e^{11})$$

261 $\langle \bar{\mathbf{a}} \nabla \phi_{0,0}, \nabla \phi_{1,0} \rangle_h = -\frac{1}{4} ((1-\lambda_e^2) \bar{a}_e^{11} + (1-\lambda_e^1) \bar{a}_e^{22}) - \frac{1}{4} \bar{a}_e^{12}$

261
$$\langle \bar{\mathbf{a}} \nabla \phi_{0,0}, \nabla \phi_{1,1} \rangle_h = -\frac{1}{4} \left((1 - \lambda_e^2) \bar{a}_e^{11} + (1 - \lambda_e^1) \bar{a}_e^{22} \right) - \frac{1}{2} \bar{a}_e^{12}$$

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

264 With (3.7) and the assumption $|\bar{a}_e^{12}| \le \min\{\bar{a}_e^{11}, \bar{a}_e^{22}\}$, we have (3.8)

$$\langle \bar{\mathbf{a}} \nabla \phi_{0,0}, \nabla \phi_{0,1} \rangle_{h} = \langle \bar{\mathbf{a}} \nabla \phi_{1,1}, \nabla \phi_{1,0} \rangle_{h} \in \left[\frac{1}{2} \left(|\bar{a}_{e}^{12}| - \bar{a}_{e}^{22} \right), \frac{1}{4} (\bar{a}_{e}^{11} - \bar{a}_{e}^{22} - |\bar{a}_{e}^{11} - \bar{a}_{e}^{22} | \right), \\ \langle \bar{\mathbf{a}} \nabla \phi_{0,0}, \nabla \phi_{1,0} \rangle_{h} = \langle \bar{\mathbf{a}} \nabla \phi_{0,1}, \nabla \phi_{1,1} \rangle_{h} \in \left[\frac{1}{2} \left(|\bar{a}_{e}^{12}| - \bar{a}_{e}^{11} \right), \frac{1}{4} (\bar{a}_{e}^{22} - \bar{a}_{e}^{11} - |\bar{a}_{e}^{11} - \bar{a}_{e}^{22} | \right), \\ \langle \bar{\mathbf{a}} \nabla \phi_{0,0}, \nabla \phi_{1,1} \rangle_{h} \in \left(-\frac{1}{2} (\min\{\bar{a}_{e}^{11}, \bar{a}_{e}^{22}\} - \bar{a}_{e}^{12}), -\frac{1}{2} (|\bar{a}_{e}^{12}| + \bar{a}_{e}^{12}) \right], \\ \langle \bar{\mathbf{a}} \nabla \phi_{0,1}, \nabla \phi_{1,0} \rangle_{h} \in \left(-\frac{1}{2} (\min\{\bar{a}_{e}^{11}, \bar{a}_{e}^{22}\} + \bar{a}_{e}^{12}), -\frac{1}{2} (|\bar{a}_{e}^{12}| - \bar{a}_{e}^{12}) \right],$$

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 values as the bound of the closed side of the interval.

Given $i, j \in \{1, \dots, N_h + N_h^\partial\}$, obviously, if both \mathbf{x}_i and \mathbf{x}_j are vertices of the

269 same elements e, then we have

$$a_{ij} = \mathcal{A}_{h}(\varphi_{j},\varphi_{i})$$

$$= \sum_{e \in \Omega_{h}} \int_{e} \bar{\mathbf{a}} \nabla \varphi_{j} \cdot \nabla \varphi_{i} d^{h}_{\boldsymbol{\lambda}_{e}} \mathbf{x} + \int_{e} c \varphi_{j} \varphi_{i} d^{h}_{1} \mathbf{x}$$

$$= \sum_{e \in \Omega_{h}} \int_{\hat{K}} \bar{\mathbf{a}} \hat{\nabla} \hat{\varphi}_{j} \cdot \hat{\nabla} \hat{\varphi}_{i} d^{h}_{\boldsymbol{\lambda}_{e}} \hat{\mathbf{x}} + \int_{\hat{K}} \hat{c} \hat{\varphi}_{j} \hat{\varphi}_{i} d^{h}_{1} \hat{\mathbf{x}}$$

$$= \sum_{i,j \in e} \int_{\hat{K}} \bar{\mathbf{a}} \hat{\nabla} \hat{\varphi}_{j} \cdot \hat{\nabla} \hat{\varphi}_{i} d^{h}_{\boldsymbol{\lambda}_{e}} \hat{\mathbf{x}} + \int_{\hat{K}} \hat{c} \hat{\varphi}_{j} \hat{\varphi}_{i} d^{h}_{1} \hat{\mathbf{x}}$$

where $\sum_{i,j\in e}$ means summation over all elements e containing both vertices i and j. 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 of the values in (3.8) depending on their relative positions. Therefore, for $i \neq j$, with (3.7) and the assumption $|\bar{a}_e^{12}| \leq \min\{\bar{a}_e^{11}, \bar{a}_e^{22}\}$ we have

275 (3.10)
$$A_{ij} = \sum_{i,j \in e} \int_{\hat{K}} \bar{\mathbf{a}} \hat{\nabla} \hat{\varphi}_j \cdot \hat{\nabla} \hat{\varphi}_i d^h_{\boldsymbol{\lambda}_e} \hat{\mathbf{x}} \le 0.$$

For $i = 1, \ldots, N_h$, we note that

277 (3.11)
$$\sum_{j=1}^{N_h+N_h^{\partial}} A_{ij} = \sum_{j=1}^{N_h+N_h^{\partial}} \mathcal{A}_h(\varphi_j,\varphi_i) = \mathcal{A}_h(1,\varphi_i) = Cc_i \ge 0,$$

where *C* is a certain positive number and $c_i = c(\mathbf{x}_i) \ge 0$. If \mathbf{x}_i has no neighboring node on the boundary, by $\sum_{j=1}^{N_h} A_{ij} = \sum_{j=1}^{N_h+N_h^{\partial}} A_{ij}$ and (3.11), the *i*-th row sum of *A* is non-negative. Therefore, we have

281 (3.12)
$$A_{ii} \ge \sum_{j=1, j \ne i}^{N_h} |A_{ij}|.$$

If \mathbf{x}_i has a neighboring node on the boundary, with (3.11) and $\sum_{j=1}^{N_h} A_{ij} \ge \sum_{j=1}^{N_h+N_h^{\partial}} A_{ij}$ due to $A_{ij} \le 0$ for $i \ne j$, we do have (3.12) holds. When \mathbf{x}_i has two neighboring node on the boundary, based on (3.8), among the two neighboring nodes on the boundary of \mathbf{x}_i , there exists nodes \mathbf{x}_l with $l \in \{N_h + 1, \dots, N_h + N_h^{\partial}\}$ such that $A_{il} < 0$. Then we have

287
$$\sum_{j=1}^{N_h} A_{ij} \ge \sum_{j=1}^{N_h + N_h^{\partial}} A_{ij} - A_{il} > 0.$$

Therefore, the stiffness matrix A is an M-matrix. Since the lumped mass matrix is diagonal and entry-wise positive, with $A_{ij} \leq 0$ and noticing that $\bar{L}_h^{-1} = \begin{pmatrix} L_h^{-1} & -L_h^{-1}B^{\partial} \\ 0 & I \end{pmatrix}$, we conclude \bar{L}_h is also an M-matrix. Then with (3.11) and Theorem 3.1 we obtain the finite difference operator defined by \bar{L}_h satisfies discrete maximum principle. 293 REMARK 1. For each element e, the choice in (3.7) make $\lambda_e^1, \lambda_e^2 > 0$, which im-294 plies the V^h -ellipticity of the bilinear form (3.3) discussed in Section 4.2. Therefore, 295 we can assure of V^h -ellipticity and the stiffness matrix being an M-matrix simultane-296 ously.

297 REMARK 2. The constraint on the coefficient, $|\bar{a}_e^{12}| \leq \min\{\bar{a}_e^{11}, \bar{a}_e^{22}\}$, aligns with 298 the condition for rendering the stiffness matrix as an *M*-matrix in the seven-point 299 stencil control volume method with optimal optimal monotonicity region in the case 300 of homogeneous medium and uniform mesh in [29]. In [28], the authors show that 301 a three-by-three stencil can be used to construct monotone finite difference schemes 302 under the assumption $|a^{12}| < \min\{a^{11}, a^{22}\}$.

REMARK 3. If the domain is not convex, e.g., an L-shaped domain, as long as it can be partitioned by uniform square meshes satisfying the coefficients constraints or quadrilateral meshes satisfying the mesh constraints derived in Section 5, the stiffness matrix is still an M-matrix and the monotonicity holds. But the a priori error estimates in Section 4 might no longer hold due to possible loss of the elliptic regularity on a nonconvex domain.

REMARK 4. The choice of the quadrature parameters in (3.7) is sharp to enforce the \bar{L}_h being an *M*-matrix but not for monotonicity since *M*-matrix property is just a sufficient but not necessary condition for monotonicity.

4. Convergence of the Q^1 finite element method with mixed quadrature. In this section, we prove the second-order accuracy of the scheme (3.2) on uniform rectangular mesh. For simplicity we only prove result for the problem with homogeneous Dirichlet boundary condition, i.e. g = 0. For convenience, in this section, we may drop the subscript h in a test function $v_h \in V^h$. When there is no confusion, we may also drop $d\mathbf{x}$ or $d\hat{\mathbf{x}}$ in a integral.

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

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

322 (4.1)
$$\int_{e} (\mathbf{a}\nabla u) \cdot \nabla v d\mathbf{x} - \int_{e} (\bar{\mathbf{a}}_{e}\nabla u) \cdot \nabla v d^{h}_{\boldsymbol{\lambda}_{e}} \mathbf{x} = \mathcal{O}(h^{2}) \|u\|_{3,e} \|v\|_{2,e},$$

323 (4.2)
$$\int_{e} cuv d\mathbf{x} - \int_{e} cuv d_{1}^{h} \mathbf{x} = \mathcal{O}\left(h^{2}\right) \|u\|_{2,e} \|v\|_{2,e}.$$

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

$$\int_{e} a u_{x_{k}} v_{x_{l}} d\mathbf{x} - \int_{e} \bar{a}_{e} u_{x_{k}} v_{x_{l}} d^{h}_{\mathbf{\lambda}_{e}} \mathbf{x}$$

$$326 \quad (4.3) \qquad = \int_{e} (a - \bar{a}_{e}) u_{x_{k}} v_{x_{l}} d\mathbf{x} + \bar{a}_{e} \left(\int_{e} u_{x_{k}} v_{x_{l}} d\mathbf{x} - \int_{e} u_{x_{k}} v_{x_{l}} d^{h}_{\mathbf{\lambda}_{e}} \mathbf{x} \right)$$

$$= \int_{e} (a - \bar{a}_{e}) u_{x_{k}} v_{x_{l}} d\mathbf{x} + \bar{a}_{e} E(u_{x_{k}} v_{x_{l}}).$$

327 For the first term,

$$\int_{e} (a - \bar{a}_{e}) u_{x_{k}} v_{x_{l}} d\mathbf{x}$$
328 (4.4)
$$= \int_{e} (a - \bar{a}_{e}) (u_{x_{k}} v_{x_{l}} - \overline{u_{x_{k}} v_{x_{l}}}) d\mathbf{x} + \int_{e} (a - \bar{a}_{e}) \overline{u_{x_{k}} v_{x_{l}}} d\mathbf{x}$$

$$\leq ||a - \bar{a}_{e}||_{0,\infty,e} ||u_{x_{k}} v_{x_{l}} - \overline{u_{x_{k}} v_{x_{l}}}||_{0,1,e} + \frac{1}{meas(e)} \int_{e} (a - \bar{a}_{e}) d\mathbf{x} \int_{e} u_{x_{k}} v_{x_{l}} d\mathbf{x}.$$

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

330 (4.5)
$$\begin{aligned} \|a - \bar{a}_e\|_{0,\infty,e} \|u_{x_k}v_{x_l} - \overline{u_{x_k}v_{x_l}}\|_{0,1,e} \\ = \mathcal{O}(h^2) \|a\|_{1,\infty,e} \|\nabla (u_{x_k}v_{x_l})\|_{0,1,e} = \mathcal{O}(h^2) \|u\|_{2,e} \|v\|_{2,e}. \end{aligned}$$

331 By Lemma 2.2 and Cauchy-Schwartz inequality

332 (4.6)
$$\frac{1}{meas(e)} \int_{e} (a - \bar{a}_{e}) d\mathbf{x} \int_{e} u_{x_{k}} v_{x_{l}} d\mathbf{x}$$
$$= \frac{h^{2+d}}{meas(e)} [a]_{2,\infty,e} ||u_{x_{k}}||_{0,e} ||v_{x_{l}}||_{0,e} = \mathcal{O}(h^{2}) ||u||_{1,e} ||v||_{1,e}$$

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

335 (4.7)
$$\int_{e} (a - \bar{a}_{e}) u_{x_{k}} v_{x_{l}} d\mathbf{x} = \mathcal{O}(h^{2}) \|a\|_{2,\infty,e} \|u\|_{2,e} \|v\|_{2,e}.$$

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

337 (4.8)
$$\int_{e} \bar{a}_{e} u_{x_{k}} v_{x_{l}} d\mathbf{x} - \int_{e} \bar{a}_{e} u_{x_{l}} v_{x_{l}} d^{h}_{\boldsymbol{\lambda}_{e}} \mathbf{x} = \mathcal{O}(h^{2}) \|a\|_{0,\infty,e} \|u\|_{3,e} \|v\|_{2,e},$$

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

339 (4.9)
$$\int_{e} a u_{x_{k}} v_{x_{l}} d\mathbf{x} - \int_{e} \bar{a}_{e} u_{x_{k}} v_{x_{l}} d^{h}_{\boldsymbol{\lambda}_{e}} \mathbf{x} = \mathcal{O}(h^{2}) \|a\|_{2,\infty,e} \|u\|_{3,e} \|v\|_{2,e}.$$

340 Therefore, we have

341 (4.10)
$$\int_{e} (\mathbf{a}(\mathbf{x})\nabla u) \cdot \nabla v d\mathbf{x} - \int_{e} (\bar{\mathbf{a}}(\mathbf{x}) \cdot \nabla u) \nabla v d_{\lambda_{e}}^{h} \mathbf{x} = \mathcal{O}(h^{2}) \|\mathbf{a}\|_{2,\infty,e} \|u\|_{3,e} \|v\|_{2,e}.$$

342 Similarly we have

343 (4.11)
$$\int_{e} cuv d\mathbf{x} - \int_{e} cuv d_{1}^{h} \mathbf{x} = \mathcal{O}\left(h^{2}\right) \|c\|_{2,\infty,e} \|u\|_{2,e} \|v\|_{2,e}.$$

344 We also have

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

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

Proof. By Theorem 4.1 and noticing that the third derivative of Q^1 polynomial 345 346vanish, we have

347 (4.12)
$$\int_{e} (\mathbf{a} \nabla v_h) \cdot \nabla w_h d\mathbf{x} - \int_{e} (\bar{\mathbf{a}}_e \nabla v_h) \cdot \nabla w_h d^h_{\boldsymbol{\lambda}_e} \mathbf{x} = \mathcal{O}(h^2) \|v_h\|_{2,e} \|w_h\|_{2,e},$$

$$\int_{e} cv_{h}w_{h}d\mathbf{x} - \int_{e} cv_{h}w_{h}d\mathbf{x} = \mathcal{O}\left(h^{2}\right)\|v_{h}\|_{2,e}\|w_{h}\|_{2,e}.$$

350 By applying the inverse estimate to polynomial z_h , we obtain

351 (4.14)
$$\int_{e} (\mathbf{a}\nabla v_{h}) \cdot \nabla w_{h} d\mathbf{x} - \int_{e} (\bar{\mathbf{a}}_{e} \nabla v_{h}) \cdot \nabla w_{h} d^{h}_{\boldsymbol{\lambda}_{e}} \mathbf{x} = \mathcal{O}(h) \|v_{h}\|_{2,e} \|w_{h}\|_{1,e},$$

352 (4.15)
$$\int_{e} cv_{h}w_{h}d\mathbf{x} - \int_{e} cv_{h}w_{h}d_{1}^{h}\mathbf{x} = \mathcal{O}\left(h\right) \|v_{h}\|_{2,e} \|w_{h}\|_{1,e}$$

Then by summing over all the elements we obtain prove the Lemma. 354

4.2. V^h-ellipticity and the dual problem. In order to prove the convergence 355 results of the scheme (3.2), we need A_h satisfies V^h -ellipticity: 356

357 (4.16)
$$\forall v_h \in V_0^h, \quad C \|v_h\|_1^2 \le A_h (v_h, v_h)$$

By following the proof of Lemma 5.1 in [16], we have 358

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

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

Proof. For element e, at first we map all the functions to the reference element K. Let $Z_{0,\hat{K}}$ denote the set of vertices on the reference element 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

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

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

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

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

$$C_{1} \left\| \hat{v}_{h} \right\|_{1,\hat{K}}^{2} \leq C_{1} \sum_{i=1}^{2} \int_{\hat{K}} (\hat{v}_{h})_{\hat{x}_{i}}^{2} d_{1}^{h} \hat{\mathbf{x}} \leq \int_{\hat{K}} (\bar{\mathbf{a}}_{e}^{ij} \nabla \hat{v}_{h}) \cdot \nabla \hat{v}_{h} d_{\boldsymbol{\lambda}_{e}}^{h} \hat{\mathbf{x}} + \int_{\hat{K}} \hat{c} \hat{v}_{h}^{2} d_{1}^{h} \hat{\mathbf{x}} \leq C_{2} \left\| \hat{v}_{h} \right\|_{1,\hat{K}}^{2}.$$

Mapping these back to the original element e and summing over all elements, by the 359

equivalence of two norms $|\cdot|_1$ and $||\cdot||_1$ for the space $H_0^1(\Omega) \supset V_0^h$, we obtain the 360conclusion. 361

This manuscript is for review purposes only.

- In the following part, we assume the assumption of Lemma 4.3 is fulfilled, i.e. the V^{h} -ellipticity holds.
- In order to apply the Aubin-Nitsche duality argument for establishing convergence of function values, we need certain estimates on a proper dual problem.
- Define $\theta := u u_h$ and consider the dual problem: find $w \in H_0^1(\Omega)$ satisfying

367 (4.17)
$$A^*(w,v) = (\theta,v), \quad \forall v \in H^1_0(\Omega),$$

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

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

- Although here the bilinear form we considered is symmetric i.e. $A(\cdot, \cdot) = A^*(\cdot, \cdot)$, we still use $A^*(\cdot, \cdot)$ for abstractness.
- 370 Let $w_h \in V_0^h$ be the solution to

371 (4.18)
$$A_h^*(w_h, v_h) = (\theta, v_h), \quad \forall v_h \in V_0^h$$

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

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

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 defined in (4.17), w_h be defined in (4.18). With elliptic regularity and V^h -ellipticity hold, we have

378 (4.19)
$$\begin{aligned} \|w - w_h\|_1 \leq Ch \|w\|_2 \\ \|w_h\|_2 \leq C \|\theta\|_0. \end{aligned}$$

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

$$C_{1} \|w_{h} - v_{h}\|_{1}^{2} \leq A_{h}^{*} (w_{h} - v_{h}, w_{h} - v_{h})$$

= $A^{*} (w - v_{h}, w_{h} - v_{h}) + [A_{h}^{*} (w_{h}, w_{h} - v_{h}) - A^{*} (w, w_{h} - v_{h})] + [A^{*} (v_{h}, w_{h} - v_{h}) - A_{h}^{*} (v_{h}, w_{h} - v_{h})]$
= $A^{*} (w - v_{h}, w_{h} - v_{h}) + [A (w_{h} - v_{h}, v_{h}) - A_{h} (w_{h} - v_{h}, v_{h})]$
 $\leq C \|w - v_{h}\|_{1} \|w_{h} - v_{h}\|_{1} + Ch \|v_{h}\|_{2} \|w_{h} - v_{h}\|_{1},$

379 which implies

380 (4.20)
$$\|w - w_h\|_1 \le \|w - v_h\|_1 + \|w_h - v_h\|_1 \le C \|w - v_h\|_1 + Ch \|v_h\|_2.$$

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

384 (4.21)
$$\|w - \Pi_1 w\|_1 \le Ch \|w\|_2, \|w - \Pi_1 w\|_2 \le C \|w\|_2,$$

- 385 which implies
- 386 (4.22) $\|\Pi_1 w\|_2 \le \|w\|_2 + \|w \Pi_1 w\|_2 \le C \|w\|_2.$

By setting $v_h = \Pi_1 w$, using (4.20), (4.21) and (4.22), we have 387

388 (4.23)
$$\|w - w_h\|_1 \le C \|w - \Pi_1 w\|_1 + Ch \|\Pi_1 w\|_2 \le Ch \|w\|_2.$$

By (4.21) and (4.23), we also have 389

390 (4.24)
$$\|w_h - \Pi_1 w\|_1 \le \|w - \Pi_1 w\|_1 + \|w - w_h\|_1 \le Ch \|w\|_2.$$

By the inverse estimate on the piece-wise polynomial $w_h - \Pi_1 w$, we obtain 391

392 (4.25)
$$||w_h||_2 \le ||w_h - \Pi_1 w||_2 + ||\Pi_1 w - w||_2 + ||w||_2 \le Ch^{-1} ||w_h - \Pi_1 w||_1 + C||w||_2.\Box$$

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

$$||w_h||_2 \le C ||w||_2 \le C ||\theta||_0$$

4.3. Convergence results. In this section, we initially establish the error es-393 timate for $||u - u_h||_{1,\Omega}$. Subsequently, we demonstrate that the Q^1 finite element method, as given by (3.2), achieves second-order accuracy for function values. 394395 We have the estimate of the error $||u - u_h||_{1,\Omega}$ as follows: 396

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

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

Proof. By the First Strang Lemma, 397 (4.26)

98
$$\|u - u_{h}\|_{1,\Omega} \leq C \left(\inf_{v_{h} \in V^{h}} \left\{ \|u - v_{h}\|_{1,\Omega} + \sup_{w_{h} \in V_{h}} \frac{|\mathcal{A}(v_{h}, w_{h}) - \mathcal{A}_{h}(v_{h}, w_{h})|}{\|w_{h}\|_{1,\Omega}} \right\} + \sup_{w_{h} \in V^{h}} \frac{|\langle f, w_{h} \rangle_{h} - (f, w_{h})|}{\|w_{h}\|_{1,\Omega}} \right).$$

39

400

By Lemma 4.2, we have: 399

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

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, 401

402
$$\|u - u_h\|_{1,\Omega} = \mathcal{O}(h)(\|u\|_{2,\Omega} + \|f\|_{2,\Omega})$$

403 In the following part we prove the Aubin-Nitsche Lemma up to the quadrature error for establishing convergence of function values. 404

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 $\theta = u - u_h \in H^1_0(\Omega)$, we have 405

406 (4.27)
$$\|\theta\|_{0}^{2} = (\theta, \theta) = A(\theta, w) = A(u - u_{h}, w_{h}) + A(u - u_{h}, w - w_{h})$$

For the first term (4.27), by Lemma 4.1, we have 407

$$A(u - u_h, w_h) = [A(u, w_h) - A_h(u_h, w_h)] + [A_h(u_h, w_h) - A(u_h, w_h)]$$

$$= (f, w_h) - \langle f, w_h \rangle_h + \mathcal{O}(h^2) ||u_h||_3 ||w_h||_2$$

$$= \mathcal{O}(h^2) ||f||_2 ||w_h||_1 + \mathcal{O}(h^2) ||u_h||_2 ||w_h||_2$$

$$= \mathcal{O}(h^2) (||f||_2 + ||u_h||_2) ||\theta||_0,$$

409 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 410 Lemma 4.4, we have 411

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

where Theorem 4.5 is used in the last inequality. Therefore, we have 413

414 (4.30)
$$A(u - u_h, w_h) = \mathcal{O}(h^2) (||f||_2 + ||u||_2) ||\theta||_0$$

For the second term (4.27), by continuity of the bilinear form and Lemma 4.4, we 415 416 have

417 (4.31)
$$A(u - u_h, w - w_h) \le C \|u - u_h\|_1 \|w - w_h\|_1 \le Ch \|u - u_h\|_1 \|w\|_2$$
$$\le Ch \|u - u_h\|_1 \|\theta\|_0 = \mathcal{O}(h^2) (\|f\|_2 + \|u\|_2) \|\theta\|_0.$$

Therefore, by (4.27), (4.28) and (4.31), we have 418

419 (4.32)
$$\|\theta\|_0 = \mathcal{O}(h^2) (\|f\|_2 + \|u\|_2).$$

REMARK 5. Similar convergence results for the Q^1 method on general quasi-420 uniform quadrilateral meshes can be established via the same proof procedure in this 421 section. 422

- 5. Extension to general quadrilateral meshes. For a quadrilateral element 423 e as in Fig. 2, let $\mathbf{F}_e = (F_{e1}, F_{e2})^T$ be 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 424
 - have

$$\nabla \varphi \circ \mathbf{F}_e = DF_e^{T-1} \hat{\nabla} \dot{\varphi}$$

 $\begin{array}{l} \mathbf{v} \mathbf{v} \circ \mathbf{r}_{e} = \mathcal{D}\mathbf{F}_{e}^{-1} \nabla \ddot{\varphi} \\ \text{where } \hat{\varphi} = \varphi \circ \mathbf{F}_{e}, \, \nabla = \left(\frac{\partial}{\partial x_{1}}, \frac{\partial}{\partial x_{2}}\right)^{T}, \, \hat{\nabla} = \left(\frac{\partial}{\partial \hat{x}_{1}}, \frac{\partial}{\partial \hat{x}_{2}}\right)^{T} \text{ and Jacobian matrix } DF_{e} = \\ \left(\begin{array}{c} \frac{\partial F_{e1}}{\partial \hat{x}_{1}} & \frac{\partial F_{e1}}{\partial \hat{x}_{2}} \\ \frac{\partial F_{e2}}{\partial \hat{x}_{1}} & \frac{\partial F_{e2}}{\partial \hat{x}_{2}} \end{array}\right). \\ \text{Theorem} \end{array}$ 425426 427 Therefore, we have

428 (5.1)
$$\int_{e} \mathbf{a} \nabla u_{h} \cdot \nabla v_{h} d\mathbf{x} = \int_{\hat{K}} \left(DF_{e}^{-1} \hat{\mathbf{a}} DF_{e}^{T-1} \hat{\nabla} \hat{u}_{h} \right) \cdot \hat{\nabla} \hat{v}_{h} \left| \det(DF_{e}) \right| d\hat{\mathbf{x}}.$$

In the case of regular meshes with mesh size h, the matrix $DF_e^{-1}\hat{\mathbf{a}}DF_e^{T-1} = \frac{1}{h^2}\hat{\mathbf{a}}$ and $\det(DF_e) = h^2$.

431 Approximate (5.1) by the mixed quadrature (2.5) with parameter $\boldsymbol{\lambda} = (\lambda^1, \lambda^2)$, 432 i.e.

433 (5.2)
$$\int_{e} (\mathbf{a}\nabla u_{h}) \cdot \nabla v_{h} \, \mathrm{d}\mathbf{x} \approx \int_{\hat{K}} \left(\tilde{\mathbf{a}}\hat{\nabla}\hat{u}_{h} \right) \cdot \hat{\nabla}\hat{v}_{h} d_{\boldsymbol{\lambda}}^{h} \hat{\mathbf{x}}$$

434 where $\tilde{\mathbf{a}} = \left(|\det(DF_e)| DF_e^{-1} \hat{\mathbf{a}} DF_e^{T-1} \right) \left(\frac{1}{2}, \frac{1}{2} \right)$. As in Fig. 2, denote

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

435 and

436
$$\overrightarrow{\mathbf{c}}_{i} = (c_{i}^{0}, c_{i}^{1})^{T}, i = 0, 1, 2, 3, \quad DF_{h} = DF_{e}(\frac{1}{2}, \frac{1}{2}), \quad \overline{\mathbf{a}}_{e} = \mathbf{a}|_{e} \circ \mathbf{F}_{e}(\frac{1}{2}, \frac{1}{2}),$$

437 then we have

438
$$DF_{h} = \frac{1}{2} \begin{pmatrix} c_{1}^{0} + c_{3}^{0} & c_{0}^{0} + c_{2}^{0} \\ c_{1}^{1} + c_{3}^{1} & c_{0}^{1} + c_{2}^{1} \end{pmatrix}, \quad DF_{h}^{-1} = \frac{1}{2 \det(DF_{h})} \begin{pmatrix} c_{0}^{1} + c_{2}^{1} & -c_{0}^{0} - c_{2}^{0} \\ -c_{1}^{1} - c_{3}^{1} & c_{1}^{0} + c_{3}^{0} \end{pmatrix},$$

438

441 and on element e, we have

442 (5.3)
$$\tilde{\mathbf{a}} = \det(DF_h)DF_h^{-1}\bar{\mathbf{a}}_e DF_h^{T-1} = \begin{pmatrix} \tilde{a}_e^{11} & \tilde{a}_e^{12} \\ \tilde{a}_e^{12} & \tilde{a}_e^{22} \end{pmatrix}.$$

To have the stiffness matrix an M-matrix, by Theorem 3.3, the following is a sufficient condition:

445 (5.4)
$$\left| \tilde{a}_{e}^{12} \right| \leq \min\{ \tilde{a}_{e}^{11}, \tilde{a}_{e}^{22} \},$$

446 where

$$\tilde{a}^{11} = \frac{1}{4|\det(DF_h)|} \begin{pmatrix} c_0^1 + c_2^1 & -c_0^0 - c_2^0 \end{pmatrix} \begin{pmatrix} \bar{a}^{11} & \bar{a}^{12} \\ \bar{a}^{12} & \bar{a}^{22} \end{pmatrix} \begin{pmatrix} c_0^1 + c_2^1 \\ -c_0^0 - c_2^0 \end{pmatrix} \\
= \frac{1}{4|\det(DF_h)|} \begin{pmatrix} c_0^0 + c_2^0 & c_0^1 + c_2^1 \end{pmatrix} \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} \bar{a}^{11} & \bar{a}^{12} \\ \bar{a}^{12} & \bar{a}^{22} \end{pmatrix} \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} \begin{pmatrix} c_0^0 + c_2^0 \\ c_0^1 + c_2^1 \end{pmatrix} \\
= \frac{1}{4|\det(DF_h)|} \begin{pmatrix} c_0^0 + c_2^0 & c_0^1 + c_2^1 \end{pmatrix} \begin{pmatrix} \bar{a}^{22} & -\bar{a}^{12} \\ -\bar{a}^{12} & \bar{a}^{11} \end{pmatrix} \begin{pmatrix} c_0^0 + c_2^0 \\ c_0^1 + c_2^1 \end{pmatrix} \\
= \frac{\det(\bar{\mathbf{a}}_e)}{4|\det(DF_h)|} (\bar{\mathbf{c}}_0^0 + \bar{\mathbf{c}}_2)^T \bar{\mathbf{a}}_e^{-1} (\bar{\mathbf{c}}_0^1 + \bar{\mathbf{c}}_2),$$

448 and similarly

449
$$\tilde{a}^{12} = -\frac{\det(\bar{\mathbf{a}}_e)}{4|\det(DF_h)|} \left(\vec{\mathbf{c}}_0 + \vec{\mathbf{c}}_2\right)^T \bar{\mathbf{a}}_e^{-1} \left(\vec{\mathbf{c}}_1 + \vec{\mathbf{c}}_3\right),$$

$$\tilde{a}^{22} = \frac{\det(\bar{\mathbf{a}}_e)}{4|\det(DF_h)|} \left(\overline{\mathbf{c}}_1^{-1} + \overline{\mathbf{c}}_3^{-1}\right)^T \bar{\mathbf{a}}_e^{-1} \left(\overline{\mathbf{c}}_1^{-1} + \overline{\mathbf{c}}_3^{-1}\right)$$

452 By
$$\overrightarrow{\mathbf{c}_1} + \overrightarrow{\mathbf{c}_2} - \overrightarrow{\mathbf{c}_3} - \overrightarrow{\mathbf{c}_0} = \overrightarrow{0}$$
, we note (5.4) is equivalent to

453 (5.5)
$$(\overrightarrow{\mathbf{c}_{0}} + \overrightarrow{\mathbf{c}_{2}})^{T} \, \overrightarrow{\mathbf{a}}_{e}^{-1} (\overrightarrow{\mathbf{c}_{0}} + \overrightarrow{\mathbf{c}_{3}}) \ge 0, \quad (\overrightarrow{\mathbf{c}_{0}} + \overrightarrow{\mathbf{c}_{2}})^{T} \, \overrightarrow{\mathbf{a}}_{e}^{-1} (\overrightarrow{\mathbf{c}_{0}} - \overrightarrow{\mathbf{c}_{1}}) \ge 0, \\ (\overrightarrow{\mathbf{c}_{1}} + \overrightarrow{\mathbf{c}_{3}})^{T} \, \overrightarrow{\mathbf{a}}_{e}^{-1} (\overrightarrow{\mathbf{c}_{0}} + \overrightarrow{\mathbf{c}_{3}}) \ge 0, \quad (\overrightarrow{\mathbf{c}_{1}} + \overrightarrow{\mathbf{c}_{3}})^{T} \, \overrightarrow{\mathbf{a}}_{e}^{-1} (\overrightarrow{\mathbf{c}_{1}} - \overrightarrow{\mathbf{c}_{0}}) \ge 0.$$



FIG. 2. A quadrilateral element e.

THEOREM 5.1. With $\tilde{\mathbf{a}}$ defined in (5.3), if the quadrilateral mesh fulfill the condi-454tion (5.4) or the mesh condition (5.5), then the stiffness matrix of the linear Q^1 finite 455element scheme (3.2) for solving BVP (1.1) is an M-matrix. 456

REMARK 6. If the diffusion coefficient matrix degenerate to a scalar, i.e. $\mathbf{a} =$ 457 $\alpha(\mathbf{x})I$, a sufficient condition for (5.5) is that both diagonals of the quadrilateral ele-458 ment bisect each angle, resulting in two non-obtuse angles for each vertex. 459

REMARK 7. By adopting some anisotropic mesh adaptation strategy where an 460 anisotropic mesh is generated as an M-uniform mesh or a uniform mesh in the metric 461 specified by the diffusion matrix \mathbf{a} . The method (3.2) for any anistropic problem 462possibly can be monotone on that anisotropic mesh. 463

If we consider rectangular meshes, for simplicity we assume 464

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

Then we have 466

465

467
$$\tilde{\mathbf{a}} = \begin{pmatrix} \frac{h_2}{h_1} \bar{a}^{11} & \bar{a}^{12} \\ \bar{a}^{12} & \frac{h_1}{h_2} \bar{a}^{22} \end{pmatrix}$$

and (5.4) becomes 468

469 (5.6)
$$|\bar{a}_e^{12}| \le \min\{\frac{h_2}{h_1}\bar{a}_e^{11}, \frac{h_1}{h_2}\bar{a}_e^{22}\}.$$

Recall that $\sqrt{\bar{a}_e^{11}\bar{a}_e^{22}} \ge |\bar{a}_e^{12}|$, taking 470

471 (5.7)
$$\frac{h_1}{h_2} = \sqrt{\frac{\bar{a}_e^{11}}{\bar{a}_e^{22}}}$$

will guarantee (5.6). Therefore, if the rectangular mesh is deployed with aspect ratio 472 $\sqrt{\frac{\bar{a}^{11}}{\bar{a}_{e}^{52}}}$, then the stiffness matrix of the Q^1 method (3.2) is an *M*-matrix. 473

If the elliptic coefficient **a** is constant on the whole domain Ω , when the rect-474 angular mesh are fine enough, there must exist rectangular mesh with aspect ratio 475approximately $\sqrt{\frac{\bar{a}_e^{11}}{\bar{a}_e^{22}}}$ such that the stiffness matrix of scheme (3.2) solve the BVP (1.1) 476is an *M*-matrix. 477

H. LI AND X. ZHANG

REMARK 8. Unfortunately, the technique in this paper cannot be easily extended 478 to the three-dimensional case. For the three-dimensional case, with the basis on the 479

reference element $\hat{K} = [0, 1]^3$ 480

481 (5.8)
$$\hat{\phi}_{i,j,k} = \hat{x}_1^i (1 - \hat{x}_1)^{1-i} \hat{x}_2^j (1 - \hat{x}_2)^{1-j} \hat{x}_3^k (1 - \hat{x}_3)^{1-k}, \, i, j, k = 0, 1,$$

and the same derivation as in two-dimensional case, we find out 482

$$483 \quad \langle \bar{\mathbf{a}} \nabla \phi_{0,0,0}, \nabla \phi_{1,1,0} \rangle_h = -\frac{1}{16} (1+\lambda_e^3) \left[(1-\lambda_e^2) \bar{a}_e^{11} + (1-\lambda_e^1) \bar{a}_e^{22} + 2\bar{a}_e^{12} \right] + \frac{1}{16} (1-\lambda_e^1) (1-\lambda_e^2) \bar{a}_e^{33}.$$

485For the symmetric positive-definite coefficient matrix

$$\mathbf{a} = \begin{pmatrix} 1 & -1+\epsilon & \epsilon \\ -1+\epsilon & 1 & -1+\epsilon \\ \epsilon & -1+\epsilon & 1 \end{pmatrix}$$

with $\frac{1}{4}(5-\sqrt{17}) < \epsilon < \frac{1}{2}$, we obtain 488

$${}^{489}_{490} \quad \langle \bar{\mathbf{a}} \nabla \phi_{0,0,0}, \nabla \phi_{1,1,0} \rangle_h = \frac{1}{16} (1 + \lambda_e^3) (\lambda_e^1 + \lambda_e^2 - 2\epsilon) + \frac{1}{16} (1 - \lambda_e^1) (1 - \lambda_e^2) \ge \frac{1}{16} (1 - 2\epsilon).$$

Then obviously the stiffness matrix fails to be an M-matrix. 491

6. Numerical experiment. 492

6.1. Numerical experiments on uniform meshes. In this subsection, we 493 show tests verifying the proved order of accuracy and monotonicity of the scheme 494(3.2) on uniform rectangular meshes. We consider the following two-dimensional 495elliptic equation with Dirichlet boundary condition: 496

497 (6.1)
$$-\nabla \cdot (\mathbf{a}\nabla u) + cu = f \quad \text{on } [0,\pi]^2$$

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

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

When solving this problem with our method, we take the quadrature parameter 498 in element e as $\lambda_e^1 = \lambda_e^2 = 1 - \frac{2|a_e^{12}|}{a_e^{11} + a_e^{22}}$. The errors are reported in Table 1. We observe the desired second-order conver-499

500gence in the discrete l^2 -norm and infinity norm for the function values. 501

The monotonicity is verified by the smallest entries in L_h^{-1} and \bar{L}_h^{-1} which are listed in Table 2. As we can see, $L_h^{-1} \ge 0$ and $\bar{L}_h^{-1} \ge 0$ are achieved. 502 503

Then we consider a more anisotropic case in the form of (6.1) with anisotropic-coef 504

505 (6.2)
$$a^{11} = 1, \quad a^{12} = a_{21} = 9.99, \quad a^{22} = 100, \quad c = x_1^2 x_2^2$$

and exact solution

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

As stated in (5.7), we set $\frac{h_1}{h_2} = \sqrt{\frac{a^{11}}{a^{22}}} = 10$, then we examine the accuracy and 506 monotonicity of the method. When solving this problem with our method, we take the quadrature parameter in element e as $\lambda_e^1 = \lambda_e^2 = 1 - \frac{2|\tilde{a}_e^{12}|}{\tilde{a}_e^{11} + \tilde{a}_e^{22}}$. 507508

The errors are reported in Table 3. We observe the desired second-order conver-509gence in the discrete l^2 -norm and infinity norm for the function values. The monotonicity is verified by the smallest entries in L_h^{-1} and \bar{L}_h^{-1} which are 510

511512listed in Table 4.

TABLE 1

A two-dimensional elliptic equation with Dirichlet boundary conditions on uniform meshes. 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}	3.56E-1	-	2.70E-2	-
8×8	7^{2}	6.41E-2	2.47	4.89E-2	2.47
16×16	15^{2}	1.49E-3	2.11	1.15E-2	2.08
32×32	31^2	3.65E-3	2.03	2.91E-3	1.99
64×64	63^{2}	9.08E-4	2.01	7.25E-4	2.00

TABLE 2 Minimum of entries in \bar{L}_h^{-1} and L_h^{-1} for elliptic equation (6.1) with smooth coefficients on uniform meshes.

FEM Mesh	\bar{L}_h^{-1}	L_h^{-1}
4×4	0	6.38E-06
8×8	0	4.26E-10
16×16	0	2.40E-14
32×32	0	1.42E-18
64×64	0	9.24E-23

6.2. Numerical experiments on quadrilateral meshes. In this subsection, 513 we show tests verifying the proved order of accuracy and monotonicity of the scheme 514(3.2) on general quadrilateral meshes. We consider the following two-dimensional 515

Poisson equation with Dirichlet boundary condition: 516

517 (6.3)
$$-\nabla \cdot (a\nabla u) + cu = f \text{ on } [0,\pi]^2$$

where $a = 1 + 10x_2^2 + x_1 \cos x_2 + x_2$ and $c = x_1^2 x_2^2$, with an exact solution

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

The domain $[0,\pi]^2$ is partitioned into $N_y \times N_x$ elements, where the elements 518are forced to adapt to an inner edge. The angle between the inner edge and the 519 *x*-axis is $\arctan(\frac{6\sqrt{3}}{5})$ as depicted in Figure 3, where $N_y = N_x = 16$. When solving this problem with our method, we take the quadrature parameter in element *e* as $\lambda_e^1 = \lambda_e^2 = 1 - \frac{2|\tilde{a}_e^{12}|}{\tilde{a}_e^{11} + \tilde{a}_e^{22}}$. The errors are reported in Table 5. We observe the desired second-order conver-520 521

523gence in the discrete l^2 -norm and infinity norm for the function values. 524

For the quadrilateral meshes in Figure 3, we can verify that (5.4) are satisfied 525 on each elements numerically. Then we verify the monotonicity through the smallest entries in L_h^{-1} and \bar{L}_h^{-1} which are listed in Table 6. As we can see, $L_h^{-1} \ge 0$ and $\bar{L}_h^{-1} \ge 0$ are achieved. 526 527 528

7. Conclusion. We constructed a linear monotone Q^1 finite element method 529 for anistropic diffusion problem (1.1). On uniform meshes, when the diffusion matrix 530 is diagonally dominant, the *M*-matrix property is guaranteed thus monotonicity is 531 achieved. When this Q^1 finite element method is deployed on a general quadrilateral 532533 mesh, we obtain a local mesh constraint.

H. LI AND X. ZHANG

TABLE 3

A two-dimensional elliptic equation with anisotropic coefficients (6.2) and Dirichlet boundary conditions on anisotropic meshes.

FEM Mesh	DoF	l^2 error	order	l^{∞} error	order
40×4	39×3	1.58E-1	-	1.20E-1	-
80×8	79×7	3.59E-2	2.14	2.72E-2	2.14
160×16	159×15	8.76E-3	2.03	6.65E-3	2.03
320×32	319×31	2.18E-3	2.01	1.65E-3	2.01
640×64	639×63	5.44E-4	2.00	4.13E-4	2.00

TABLE 4

Minimum of entries in \bar{L}_h^{-1} and L_h^{-1} for elliptic equation (6.1) with anisotropic coefficients (6.2) on anisotropic meshes.

FEM Mesh	\bar{L}_h^{-1}	L_h^{-1}
40×4	0	0
80×8	0	0
160×16	0	0
320×32	0	0
640×64	0	0

534

REFERENCES

- [1] J. BRAMBLE, B. HUBBARD, AND V. THOMÉE, Convergence estimates for essentially positive 536type discrete dirichlet problems, Mathematics of Computation, 23 (1969), pp. 695–709.
- 537[2] C. CANCÈS, M. CATHALA, AND C. LE POTIER, Monotone corrections for generic cell-centered 538 finite volume approximations of anisotropic diffusion equations, Numerische Mathematik, 539125 (2013), pp. 387-417.
- 540[3] I. CHRISTIE AND C. HALL, The maximum principle for bilinear elements, Internat. J. Numer. 541Methods Engrg., 20 (1984), pp. 549-553.
- 542[4] P. G. CIARLET, Discrete maximum principle for finite-difference operators, Aequationes Math., 5434 (1970), pp. 338-352.
- 544[5] P. G. CIARLET, The finite element method for elliptic problems, Classics in applied mathemat-545ics, 40 (2002), pp. 1–511.
- 546[6] L. J. CROSS AND X. ZHANG, On the monotonicity of Q^2 spectral element method for laplacian on quasi-uniform rectangular meshes, Communications in Computational Physics, 35 (2024), 547548pp. 160–180.
- [7] L. J. CROSS AND X. ZHANG, On the monotonicity of Q^3 spectral element method for laplacian, 549550Annals of Applied Mathematics, 40 (2024), pp. 161–190.
- 551[8] L. C. EVANS, Partial Differential Equations, vol. 019 of Graduate Studies in Mathematics, 552American Mathematical Society, 2 ed., 2010.
- [9] D. GREENSPAN AND P. JAIN, On non negative difference analogues of elliptic differential equa-553554tions, Journal of the Franklin Institute, 279 (1965), pp. 360–365.
- 555[10] J. HU AND X. ZHANG, Positivity-preserving and energy-dissipative finite difference schemes 556for the Fokker-Planck and Keller-Segel equations, IMA Journal of Numerical Analysis, 43 557(2022), pp. 1450–1484.
- [11] W. HUANG, Discrete maximum principle and a delaunay-type mesh condition for linear fi-558559nite element approximations of two-dimensional anisotropic diffusion problems, Numerical 560Mathematics: Theory, Methods and Applications, 4 (2011), pp. 319-334.
- 561[12] D. KUZMIN, M. J. SHASHKOV, AND D. SVYATSKIY, A constrained finite element method satisfying the discrete maximum principle for anisotropic diffusion problems, J. Comput. Phys., 562563 228 (2009), pp. 3448-3463.
- 564[13] C. LE POTIER, A nonlinear finite volume scheme satisfying maximum and minimum principles 565for diffusion operators, International Journal on Finite Volumes, (2009), pp. 1–20.
- [14] H. LI, S. XIE, AND X. ZHANG, A high order accurate bound-preserving compact finite difference 566567 scheme for scalar convection diffusion equations, SIAM Journal on Numerical Analysis,



FIG. 3. Quadrilateral mesh.

 $\begin{array}{c} {\rm TABLE \ 5} \\ A \ two-dimensional \ Poisson \ equation \ with \ Dirichlet \ boundary \ conditions \ on \ quadrilateral \ meshes \end{array}$

FEM Mesh	DoF	l^2 error	order	l^{∞} error	order
4×4	3^{2}	1.24E-1	-	8.70E-2	-
8×8	7^{2}	3.19E-2	1.96	2.84E-2	1.61
16×16	15^{2}	7.82E-3	2.03	6.93E-3	2.04
32×32	31^{2}	1.94E-3	2.01	1.76E-3	1.97
64×64	63^{2}	4.85E-4	2.00	4.41E-4	2.00

- 568 56 (2018), pp. 3308–3345.
- 569 [15] H. LI AND X. ZHANG, On the monotonicity and discrete maximum principle of the finite 570 difference implementation of C^0 - Q^2 finite element method, Numerische Mathematik, 145 571 (2020), pp. 437–472.
- [16] H. LI AND X. ZHANG, Superconvergence of high order finite difference schemes based on variational formulation for elliptic equations, Journal of Scientific Computing, 82 (2) (2020).
- [17] H. LI AND X. ZHANG, A high order accurate bound-preserving compact finite difference scheme
 for two-dimensional incompressible flow, Communications on Applied Mathematics and
 Computation, (2023), pp. 1–29.
- [18] X. LI AND W. HUANG, An anisotropic mesh adaptation method for the finite element solution of heterogeneous anisotropic diffusion problems, Journal of Computational Physics, 229 (2010), pp. 8072–8094.
- [19] X. LI, D. SVYATSKIY, AND M. SHASHKOV, Mesh adaptation and discrete maximum principle
 for 2d anisotropic diffusion problems, tech. report, Technical Report LA-UR 10-01227, Los
 Alamos National Laboratory, Los Alamos, NM, 2007.
- [20] K. LIPNIKOV, M. SHASHKOV, D. SVYATSKIY, AND Y. VASSILEVSKI, Monotone finite volume schemes for diffusion equations on unstructured triangular and shape-regular polygonal meshes, Journal of Computational Physics, 227 (2007), pp. 492–512.
- [21] R. LISKA AND M. SHASHKOV, Enforcing the discrete maximum principle for linear finite element solutions of second-order elliptic problems, Commun. Comput. Phys., 3 (2008), pp. 852– 877.
- [22] C. LIU, Y. GAO, AND X. ZHANG, Structure preserving schemes for fokker-planck equations of irreversible processes, Journal of Scientific Computing, 98 (4) (2024).

TABLE 6

Minimum of entries in \bar{L}_h^{-1} and L_h^{-1} for elliptic equation (6.1) on quadrilateral meshes.

FEM Mesh	\bar{L}_h^{-1}	L_h^{-1}
4×4	0	2.31E-4
8×8	0	1.47E-5
16×16	0	8.06E-7
32×32	0	4.64E-8
64×64	0	2.77E-9

- [23] C. LIU AND X. ZHANG, A positivity-preserving implicit-explicit scheme with high order poly nomial basis for compressible Navier-Stokes equations, Journal of Computational Physics,
 493 (2023), p. 112496.
- 594 [24] J. LORENZ, Zur inversionotonie diskreter probleme, Numer. Math., 27 (1977), pp. 227–238.
- [25] C. LU, W. HUANG, AND J. QIU, Maximum principle in linear finite element approximations of anisotropic diffusion-convection-reaction problems, Numerische Mathematik, 127 (2014), pp. 515-537.
- [26] M. J. MLACNIK AND L. J. DURLOFSKY, Unstructured grid optimization for improved monotonicity of discrete solutions of elliptic equations with highly anisotropic coefficients, Journal of Computational Physics, 216 (2006), pp. 337–361.
- [27] T. S. MOTZKIN AND W. WASOW, On the approximation of linear elliptic differential equations
 by difference equations with positive coefficients, Journal of Mathematics and Physics, 31
 (1952), pp. 253–259.
- [28] C. NGO AND W. HUANG, Monotone finite difference schemes for anisotropic diffusion problems via nonnegative directional splittings, Communications in Computational Physics, 19 (2016), pp. 473–495.
- [29] J. M. NORDBOTTEN, I. AAVATSMARK, AND G. EIGESTAD, Monotonicity of control volume methods, Numerische Mathematik, 106 (2007), pp. 255–288.
- [30] R. J. PLEMMONS, M-matrix characterizations. I-nonsingular M-matrices, Numer. Anal. Appl.,
 18 (1977), pp. 175–188.
- [31] P. SHARMA AND G. W. HAMMETT, Preserving monotonicity in anisotropic diffusion, Journal
 of Computational Physics, 227 (2007), pp. 123–142.
- [32] J. SHEN AND X. ZHANG, Discrete maximum principle of a high order finite difference scheme for
 a generalized Allen-Cahn equation, Communications in Mathematical Sciences, 20 (2022),
 pp. 1409–1436.
- [33] J. WEICKERT ET AL., Anisotropic diffusion in image processing, vol. 1, Teubner Stuttgart, 1998.
- [34] J. XU AND L. ZIKATANOV, A monotone finite element scheme for convection-diffusion equations, Math. Comp., 68 (1999), pp. 1429–1446.
- [35] G. YUAN AND Z. SHENG, Monotone finite volume schemes for diffusion equations on polygonal
 meshes, Journal of computational physics, 227 (2008), pp. 6288–6312.