

C^0 PENALTY METHODS FOR THE FULLY NONLINEAR MONGE-AMPÈRE EQUATION

SUSANNE C. BRENNER, THIRUPATHI GUDI, MICHAEL NEILAN, AND LI-YENG SUNG

ABSTRACT. In this paper, we develop and analyze C^0 penalty methods for the fully nonlinear Monge-Ampère equation $\det(D^2u) = f$ in two dimensions. The key idea in designing our methods is to build discretizations such that the resulting discrete linearizations are symmetric, stable, and consistent with the continuous linearization. We are then able to show the well-posedness of the penalty method as well as quasi-optimal error estimates using the Banach fixed-point theorem as our main tool. Numerical experiments are presented which support the theoretical results.

1. INTRODUCTION

Consider the following boundary-value problem for the Monge-Ampère equation: [25, 23, 10]:

$$(1.1a) \quad \det(D^2u) = f \quad \text{in } \Omega,$$

$$(1.1b) \quad u = g \quad \text{on } \partial\Omega.$$

Here, $\Omega \subset \mathbf{R}^2$ is a convex domain, f is a strictly positive function on Ω , and

$$\det(D^2u) = \frac{\partial^2 u}{\partial x_1^2} \frac{\partial^2 u}{\partial x_2^2} - \left(\frac{\partial^2 u}{\partial x_1 \partial x_2} \right)^2$$

denotes the determinant of the Hessian matrix D^2u .

The goal of this paper is to develop and analyze C^0 penalty methods for classical solutions of the Monge-Ampère equation. In particular, we assume that (1.1) has a strictly convex solution $u \in H^s(\Omega)$ with $s > 3$. We will also assume that Ω is either a convex polygon or a smooth convex domain. In the case where Ω is smooth, the smoothness of u follows from the smoothness of f and g by the results in Caffarelli-Nirenberg-Spruck [10]. In fact, in this case there are exactly two solutions: one being convex, the other concave [14]. The case of less smooth (viscosity) solutions as well as more general types of Monge-Ampère equations (where f depends on u and ∇u) and the three-dimensional case will be addressed in future works.

Even for linear problems, the case of curved boundaries require special care in the finite element context. To cope with this difficulty, we use Nitsche's method

Received by the editor May 9, 2010 and, in revised form, September 1, 2010.

2010 *Mathematics Subject Classification*. Primary 65N30, 35J60.

Key words and phrases. Monge-Ampère equation, fully nonlinear PDEs, finite element method, convergence analysis.

The work of the first and fourth authors were supported in part by the National Science Foundation under Grants No. DMS-07-13835 and DMS-10-16332. The work of the third author was supported by the National Foundation under Grant No. DMS-09-02683.

[28] to enforce the Dirichlet boundary condition (1.1b). The crux of this method is that a weak form of the essential boundary condition (1.1b) is included into the variational formulation by penalization techniques.

Another hurdle (and the most obvious one) that must be overcome is the strong nonlinearity in the PDE (1.1a). In order to deal with this difficulty, in both the construction and analysis of the numerical scheme, we derive consistent discretizations of (1.1) such that the resulting (discrete) linearization is stable. This simple idea leads to not-so-obvious discretizations of (1.1). For example, the most natural finite element method using Nitsche's technique is to find a finite element function u_h such that

$$\int_{\Omega} (f - \det(D_h^2 u_h)) v \, dx + \sigma \sum_{e \in \mathcal{E}_h^b} h_e^{-1} \int_e (u_h - g) v \, ds = 0$$

for all test functions v . However this naive method does not work either in practice or analysis, as the linearization of the discrete problem is not consistent with the linearization of the continuous problem. In the next section, we derive some consistent schemes which inherit a stable discrete linearization.

Once we are able to construct such methods, we carry out the numerical analysis using the Banach fixed-point theorem as our main tool. We prove optimal order error estimates in H^1 -like and H^2 -like norms. Furthermore, by way of a duality argument we are able to derive (almost) optimal order error estimates in the L^2 norm.

Due to their important role in many application areas [11, 31, 9] there has been a growing interest and a surge of papers in recent years towards developing numerical schemes for the Monge-Ampère equation. The first attempt was by Oliker and Prussner [30] who constructed schemes for computing the Aleksandrov measure induced by D^2u and obtained the solution of problem (1.1) as a by-product. Using convexity arguments, they showed that their sequences converge monotonically to the solution, although no rates of convergence were provided. More recently, Oberman [29, 22] constructed a wide stencil difference scheme for nonlinear elliptic PDEs which can be written as functions of eigenvalues of the Hessian matrix (such as the Monge-Ampère equation). It was proved that the finite difference scheme satisfies the convergence criterion (consistency, stability, and monotonicity) established by Barles and Souganidis [2], although no rates of convergence were given. Other relevant finite difference schemes include [3] and [1].

The implementation and the convergence theory of finite element methods for the Monge-Ampère equation is less understood. Dean and Glowinski [15] presented an augmented Lagrange multiplier method and a least squares method for the Monge-Ampère equation by treating the nonlinear equations as a constraint and using a variational principle to select a particular solution. The convergence of their method still remains an open problem. Böhmer [5] introduced a projection method using C^1 finite element functions for a certain class of fully nonlinear second order elliptic PDEs and analyzes the methods using consistency and stability arguments. Finally, Feng and the third author considered fourth order singular perturbations of (1.1) by adding a small multiple of the biharmonic operator to the PDE (1.1a) [18]. Many numerical methods for the regularized problem were proposed in [19, 20, 27].

The finite element method considered here is a projection-type method as is the method proposed in [5]. But we use the standard Lagrange finite elements

which were originally designed for second order linear problems. Advantages of our method in comparison to the aforementioned works include:

- Lagrange elements are simple to use and are available on practically all finite element commercial software.
- The method has a similar number of degrees of freedom compared to non-conforming finite element methods [27], but is easier to implement.
- The method has a small number of degrees of freedom in comparison to C^1 finite element methods [5, 20] and mixed finite element methods [19].
- Curved boundaries can be handled easily in comparison to finite difference methods [29, 30, 3].
- We can handle general Monge-Ampère equations, where the function f depends on both u and its gradient.
- The method and analysis can be extended to the three-dimensional setting [8].

The rest of the paper is organized as follows. In Section 2, we set the notation and give the motivation behind the C^0 penalty method. We end this section with a few standard lemmas that are used frequently throughout the paper. In Section 3, we give a complete analysis of the discrete nonlinear problem. Using the Banach fixed point theorem, we are able to establish the well-posedness of the method as well as derive quasi-optimal error estimates. We end this section by deriving L^2 estimates using a duality argument. Finally, in Section 4 we provide some numerical examples which show the efficiency of the method as well as back up the theoretical findings.

2. NOTATION AND DISCRETIZATION

Throughout the paper, we use $H^r(\Omega)$ ($r \geq 0$) to denote the set of all $L^2(\Omega)$ functions whose distributional derivatives up to order r are in $L^2(\Omega)$, and we use $H_0^r(\Omega)$ to denote the set of functions whose traces vanish up to order $r - 1$ on $\partial\Omega$. For a normed linear space X , we denote by X' its dual and $\langle \cdot, \cdot \rangle$ the pairing between X and X' .

We let \mathcal{T}_h be a quasi-uniform, simplicial, and conforming triangulation [12, 7, 4] of the domain Ω where each triangle on the boundary has at most one curved side. Set $h_T = \text{diam}(T) \forall T \in \mathcal{T}_h$, $h = \max_{T \in \mathcal{T}_h} h_T$, and define the piecewise Sobolev space associated with the mesh as

$$H^r(\Omega; \mathcal{T}_h) = \prod_{T \in \mathcal{T}_h} H^r(T).$$

We denote by \mathcal{E}_h^i the set of interior edges in \mathcal{T}_h , \mathcal{E}_h^b the set of boundary edges, $\mathcal{E}_h = \mathcal{E}_h^i \cup \mathcal{E}_h^b$, and $h_e = \text{diam}(e) \forall e \in \mathcal{E}_h$.

We define the jump of a vector function \mathbf{w} on an interior edge $e = \partial T^+ \cap \partial T^-$ as follows:

$$[[\mathbf{w}]]|_e = \mathbf{w}^+ \cdot \mathbf{n}_+|_e + \mathbf{w}^- \cdot \mathbf{n}_-|_e \in \mathbf{R},$$

where $\mathbf{w}^\pm = \mathbf{w}|_{T^\pm}$ and \mathbf{n}_\pm is the outward unit normal of T^\pm . On a boundary edge $e \in \mathcal{E}_h^b$, we define

$$[[\mathbf{w}]]|_e = \mathbf{w} \cdot \mathbf{n}|_e \in \mathbf{R}.$$

Next, for a matrix $\underline{\mathbf{w}} \in \mathbf{R}^{2 \times 2}$, we define the average of $\underline{\mathbf{w}}$ on $e = \partial T^+ \cap \partial T^-$ by

$$\{\{\underline{\mathbf{w}}\}\}_e = \frac{1}{2} (\underline{\mathbf{w}}^+|_e + \underline{\mathbf{w}}^-|_e) \in \mathbf{R}^{2 \times 2}.$$

On the boundary $e \in \mathcal{E}_h^b$, we take

$$\{\{\underline{\mathbf{w}}\}\}_e = \underline{\mathbf{w}}|_e \in \mathbf{R}^{2 \times 2}.$$

We denote by $D_h^2 w$ the piecewise Hessian of a function $w \in H^2(\Omega; \mathcal{T}_h)$ and use $\text{cof}(D_h^2 w)$ to denote the piecewise cofactor matrix of $D_h^2 w$; i.e.,

$$\text{cof}(D_h^2 w)|_T = \begin{pmatrix} \frac{\partial^2 w_T}{\partial x_2^2} & -\frac{\partial^2 w_T}{\partial x_1 \partial x_2} \\ -\frac{\partial^2 w_T}{\partial x_1 \partial x_2} & \frac{\partial^2 w_T}{\partial x_1^2} \end{pmatrix} \quad \forall T \in \mathcal{T}_h,$$

where $w_T = w|_T$.

To construct a consistent scheme which inherits a stable discrete linearization, we take $w \in H^3(\Omega; \mathcal{T}_h) \cap H^1(\Omega)$ and $v \in H^2(\Omega; \mathcal{T}_h) \cap H^1(\Omega)$. By using (1.1a), the divergence-free row property of cofactor matrices (cf. Lemma 2.3), and integration by parts, we have

$$\begin{aligned} & \int_{\Omega} (f - \det(D_h^2(u + w)))v \, dx \\ &= - \int_{\Omega} (\det(D_h^2 w))v \, dx - \sum_{T \in \mathcal{T}_h} \int_T (\nabla \cdot (\text{cof}(D^2 u) \nabla w))v \, dx \\ &= - \int_{\Omega} (\det(D_h^2 w))v \, dx + \int_{\Omega} (\text{cof}(D^2 u) \nabla w) \cdot \nabla v \, dx - \sum_{e \in \mathcal{E}_h} \int_e [[\text{cof}(D^2 u) \nabla w]]v \, ds. \end{aligned}$$

Therefore, by rearranging and adding terms on both sides of the equation, we get

$$\begin{aligned} (2.1) \quad & \int_{\Omega} (f - \det(D_h^2(u + w)))v \, dx + \sum_{e \in \mathcal{E}_h^i} \int_e [[\{\{\text{cof}(D_h^2(u + w))\}\} \nabla(u + w)]]v \, ds \\ &= \int_{\Omega} (\text{cof}(D^2 u) \nabla w) \cdot \nabla v \, dx - \sum_{e \in \mathcal{E}_h^b} \int_e [[\text{cof}(D^2 u) \nabla w]]v \, ds \\ &\quad - \int_{\Omega} (\det(D_h^2 w))v \, dx + \sum_{e \in \mathcal{E}_h^i} \int_e [[\{\{\text{cof}(D_h^2 w)\}\} \nabla w]]v \, ds. \end{aligned}$$

The guiding principle here and in the derivation of (2.2) below is that we want to write the left-hand side as a functional in $u + w$ and the right-hand side as the sum of a linear and a quadratic functional in w (cf. (2.3) below).

Observe that the bilinear form

$$(w, v) \rightarrow \int_{\Omega} (\text{cof}(D^2 u) \nabla w) \cdot \nabla v \, dx - \sum_{e \in \mathcal{E}_h^b} \int_e [[\text{cof}(D^2 u) \nabla w]]v \, ds$$

that appears on the right-hand side of (2.1) can be symmetrized and stabilized to become a consistent and stable bilinear form for the second order differential operator $-\nabla \cdot (\text{cof}(D^2 u) \nabla(\cdot))$, as developed by Nitsche in [28]. After symmetrization

and stabilization, equation (2.1) becomes

$$\begin{aligned}
 (2.2) \quad & \int_{\Omega} (f - \det(D_h^2(u+w)))v \, dx + \sum_{e \in \mathcal{E}_h^i} \int_e [\{\{\operatorname{cof}(D_h^2(u+w))\}\} \nabla(u+w)]v \, ds \\
 & - \sum_{e \in \mathcal{E}_h^b} [\operatorname{cof}(D_h^2(u+w)) \nabla v](u+w) \, ds + \sum_{e \in \mathcal{E}_h^b} \int_e [\operatorname{cof}(D_h^2(u+w)) \nabla v]g \, ds \\
 & + \sigma \sum_{e \in \mathcal{E}_h^b} h_e^{-1} \int_e (u+w)v \, ds - \sigma \sum_{e \in \mathcal{E}_h^b} h_e^{-1} \int_e gv \, ds \\
 & = \int_{\Omega} (\operatorname{cof}(D^2u) \nabla w) \cdot \nabla v \, dx - \sum_{e \in \mathcal{E}_h^b} \int_e [\operatorname{cof}(D^2u) \nabla w]v \, ds \\
 & - \sum_{e \in \mathcal{E}_h^b} \int_e [\operatorname{cof}(D^2u) \nabla v]w \, ds + \sigma \sum_{e \in \mathcal{E}_h^b} h_e^{-1} \int_e vw \, ds \\
 & - \int_{\Omega} (\det(D_h^2w))v \, dx + \sum_{e \in \mathcal{E}_h^i} \int_e [\{\{\operatorname{cof}(D_h^2w)\}\} \nabla w]v \, ds \\
 & - \sum_{e \in \mathcal{E}_h^b} \int_e [\operatorname{cof}(D_h^2w) \nabla v]w \, ds,
 \end{aligned}$$

where σ is a positive penalty parameter. Equation (2.2) can be written compactly as

$$(2.3) \quad F(u+w) = Lw + Rw,$$

where the nonlinear mappings $F, R : H^3(\Omega; \mathcal{T}_h) \rightarrow [H^2(\Omega; \mathcal{T}_h) \cap H^1(\Omega)]'$ and the linear mapping $L : H^2(\Omega; \mathcal{T}_h) \rightarrow [H^2(\Omega; \mathcal{T}_h) \cap H^1(\Omega)]'$ are defined by

$$\begin{aligned}
 (2.4) \quad \langle Fw, v \rangle & = \int_{\Omega} (f - \det(D_h^2w))v \, dx + \sum_{e \in \mathcal{E}_h^i} \int_e [\{\{\operatorname{cof}(D_h^2w)\}\} \nabla w]v \, ds \\
 & - \sum_{e \in \mathcal{E}_h^b} \int_e [\operatorname{cof}(D_h^2w) \nabla v](w-g) \, ds + \sigma \sum_{e \in \mathcal{E}_h^b} h_e^{-1} \int_e (w-g)v \, ds,
 \end{aligned}$$

$$\begin{aligned}
 (2.5) \quad \langle Rw, v \rangle & = - \int_{\Omega} (\det(D_h^2w))v \, dx + \sum_{e \in \mathcal{E}_h^i} \int_e [\{\{\operatorname{cof}(D_h^2w)\}\} \nabla w]v \, ds \\
 & - \sum_{e \in \mathcal{E}_h^b} \int_e [\operatorname{cof}(D_h^2w) \nabla v]w \, ds,
 \end{aligned}$$

$$\begin{aligned}
 (2.6) \quad \langle Lw, v \rangle & = \int_{\Omega} (\operatorname{cof}(D^2u) \nabla w) \cdot \nabla v \, dx - \sum_{e \in \mathcal{E}_h^b} \int_e [\operatorname{cof}(D^2u) \nabla w]v \, ds \\
 & - \sum_{e \in \mathcal{E}_h^b} \int_e [\operatorname{cof}(D^2u) \nabla v]w \, ds + \sigma \sum_{e \in \mathcal{E}_h^b} h_e^{-1} \int_e vw \, ds.
 \end{aligned}$$

Lastly, to define the finite element method, we must define the appropriate finite element spaces. For an integer $k \geq 3$, we define the finite element space $V_h \subset H^1(\Omega)$ as follows:

- If $T \in \mathcal{T}_h$ does not have a curved edge, then $v|_T$ is a polynomial of (total) degree $\leq k$ in the rectilinear coordinates for T ;
- If $T \in \mathcal{T}_h$ has one curved edge, then $v|_T$ is a polynomial of degree $\leq k$ in the curvilinear coordinates of T that correspond to the rectilinear coordinates on the reference triangle (Example 2, p. 1216 of [4]).

Remark 2.1. The requirement $k \geq 3$ as well as the regularity condition $u \in H^s(\Omega)$, $s > 3$ will be made obvious in Theorem 3.1.

Let $F_h : V_h \rightarrow V'_h$ be the restriction of F to V_h . Then the penalty method for (1.1) is to find $u_h \in V_h$ such that

$$(2.7) \quad F_h u_h = 0;$$

that is,

$$\int_{\Omega} (f - \det(D_h^2 u_h)) v \, dx + \sum_{e \in \mathcal{E}_h^i} \int_e [\{\text{cof}(D_h^2 u_h)\} \nabla u_h] v \, ds - \sum_{e \in \mathcal{E}_h^b} \int_e [\{\text{cof}(D_h^2 u_h) \nabla v\}] (u_h - g) \, ds + \sigma \sum_{e \in \mathcal{E}_h^b} h_e^{-1} \int_e (u_h - g) v \, ds = 0 \quad \forall v \in V_h.$$

The well-posedness as well as the error estimates of the penalty method (2.7) are established in the following section. For now, we end the current section with a few remarks and technical lemmas which are used throughout the paper.

Remark 2.2. Noting R is quadratic in its arguments, we conclude from (2.3) that the operator L is the linearization of F at u . The motivation of the finite element method (2.7) is based on the fact that L is the operator associated with Nitsche’s method for the second order differential operator $-\nabla \cdot (\text{cof}(D^2 u) \nabla(\cdot))$, i.e., the linearization of the nonlinear operator (1.1a).

Remark 2.3. In order to avoid the proliferation of constants, we shall use the notation $A \lesssim B$ to represent the relation $A \leq \text{constant} \times B$, where the constant is independent of the mesh parameter h and the penalty parameter σ .

Lemma 2.1 (An Algebraic Identity). *For any $v, w \in H^2(\Omega; \mathcal{T}_h)$, there holds*

$$(2.8) \quad \det(D_h^2 v) - \det(D_h^2 w) = \frac{1}{2} \text{cof}(D_h^2 v) : D_h^2 v - \frac{1}{2} \text{cof}(D_h^2 w) : D_h^2 w = \frac{1}{2} (\text{cof}(D_h^2 v) + \text{cof}(D_h^2 w)) : (D_h^2 v - D_h^2 w).$$

Lemma 2.2 (A Discrete Sobolev Inequality [6, 7]). *For any $v \in V_h$ there holds*

$$(2.9) \quad \|v\|_{L^\infty(\Omega)} \lesssim (1 + |\ln h|^{\frac{1}{2}}) \|v\|_{H^1(\Omega)}.$$

Lemma 2.3 (Divergence-Free Property of Cofactor Matrices [17]). *For any smooth function v ,*

$$(2.10) \quad \nabla \cdot (\text{cof}(D^2 v)_i) = \sum_{j=1}^2 \frac{\partial}{\partial x_j} (\text{cof}(D^2 v)_{ij}) = 0 \quad \text{for } i = 1, 2,$$

where $\text{cof}(D^2 v)_i$ and $\text{cof}(D^2 v)_{ij}$ denote respectively the i th row and the (i, j) -entry of the cofactor matrix $\text{cof}(D^2 v)$.

Lemma 2.4 (Trace and Inverse Inequalities [7, 12]). *Let $D \subset \mathbf{R}^2$ be a regular and star-like domain. Then there holds*

$$\begin{aligned} \|v\|_{L^2(\partial D)} &\lesssim \text{diam}(D)^{\frac{1}{2}}\|v\|_{H^1(D)} + \text{diam}(D)^{-\frac{1}{2}}\|v\|_{L^2(D)} & \forall v \in H^1(D), \\ \|v\|_{W^{m,q}(D)} &\lesssim \text{diam}(D)^{n-m+2\min\{0,1/q-1/p\}}\|v\|_{W^{n,p}(D)} & \forall v \in \mathbb{P}_k(D), \end{aligned}$$

where $\mathbb{P}_k(D)$ denotes the set of polynomials up to degree k restricted to D .

Lemma 2.5 (Approximation Properties of V_h [4]). *Let m, ℓ be two integers such that $0 \leq m \leq \ell \leq k + 1$. Then for any $\chi \in H^\ell(\Omega)$, there exists $v \in V_h$ such that*

$$\left(\sum_{T \in \mathcal{T}_h} \|\chi - v\|_{H^m(T)}^2 \right)^{\frac{1}{2}} \lesssim h^{\ell-m} \|\chi\|_{H^\ell(\Omega)}.$$

Remark 2.4. Due to interpolation of Sobolev spaces (e.g. [7, Theorem 14.3.3]), the parameter ℓ appearing in Lemma 2.5 can be taken to be any real number such that $0 \leq m \leq \ell \leq k + 1$.

3. CONVERGENCE ANALYSIS

Let $L_h : V_h \rightarrow V'_h$ be the restriction of L (defined by (2.6)) to V_h , and let $L_h^{-1} : V'_h \rightarrow V_h$ denote its inverse (which exists if σ is sufficiently large; see Lemma 3.1). Define the mapping $M : H^3(\Omega; \mathcal{T}_h) \rightarrow V_h$ as

$$(3.1) \quad M = L_h^{-1}(L - F),$$

and let $M_h : V_h \rightarrow V_h$ be the restriction of M to V_h ; that is,

$$(3.2) \quad M_h = Id_h - L_h^{-1}F_h,$$

where Id_h is the identity map on V_h . The existence of a solution to (2.7) near u will be proven by establishing a fixed point for M_h in a small ball centered at $u_{c,h} \in V_h$, where

$$(3.3) \quad u_{c,h} = L_h^{-1}Lu.$$

Define the discrete H^1 -like and H^2 -like norms as

$$(3.4) \quad \|v\|_{1,h}^2 = \|v\|_{H^1(\Omega)}^2 + \sum_{e \in \mathcal{E}_h^b} \left(h_e^{-1} \|v\|_{L^2(e)}^2 + h_e \|\nabla v\|_{L^2(e)}^2 \right),$$

$$(3.5) \quad \|v\|_{2,h}^2 = \sum_{T \in \mathcal{T}_h} |v|_{H^2(T)}^2 + \sum_{e \in \mathcal{E}_h} \left(h_e^{-1} \|[\![\nabla v]\!] \|_{L^2(e)}^2 + h_e \| \{ \{ D_h^2 v \} \} \|_{L^2(e)}^2 \right) + \sum_{e \in \mathcal{E}_h^b} h_e^{-3} \|v\|_{L^2(e)}^2,$$

and define the corresponding discrete negative norm as

$$(3.6) \quad \|q\|_{-1,h} = \sup_{0 \neq v \in V_h} \frac{\langle q, v \rangle}{\|v\|_{1,h}}.$$

Remark 3.1. By Lemma 2.4 and scaling, there holds

$$(3.7) \quad \|v\|_{2,h} \lesssim h^{-1} \|v\|_{1,h} \quad \forall v \in V_h.$$

Remark 3.2. By Lemma 2.5 and scaling, if $u \in H^s(\Omega)$, then there exists $v \in V_h$ such that

$$(3.8) \quad \|u - v\|_{1,h} + h\|u - v\|_{2,h} \lesssim h^{\ell-1} \|u\|_{H^\ell(\Omega)},$$

where $\ell = \min\{k + 1, s\}$ and k denotes the polynomial degree of the finite element space V_h .

The analysis of the discrete problem (2.7) is based on the following lemmas.

Lemma 3.1. *We have*

$$(3.9) \quad \|Lv\|_{-1,h} \lesssim (1 + \sigma)\|v\|_{1,h} \quad \forall v \in H^2(\Omega; \mathcal{T}_h) \cap H^1(\Omega).$$

Moreover, there exists a $\sigma_0 > 0$ that depends on u and the minimum angle of \mathcal{T}_h such that the map $L_h : V_h \rightarrow V'_h$ is invertible and

$$(3.10) \quad \|L_h^{-1}q\|_{1,h} \lesssim \|q\|_{-1,h},$$

provided $\sigma \geq \sigma_0$.

Proof. Since u is strictly convex in Ω , the matrix $\text{cof}(D^2u)$ is positive definite. Furthermore by a Sobolev embedding, if $u \in H^s(\Omega)$ for some $s > 3$, then $u \in W^{2,\infty}(\Omega)$. Therefore, there exist constants $\lambda, \Lambda > 0$ such that

$$(3.11) \quad \lambda \|w\|_{H^1(\Omega)}^2 \leq \int_{\Omega} (\text{cof}(D^2u)\nabla w) \cdot \nabla w \, dx \leq \Lambda \|w\|_{H^1(\Omega)}^2 \quad \forall w \in H^1(\Omega).$$

The conclusion of the lemma then follows from (3.11) and the results of Nitsche [28]. □

Remark 3.3. By Lemma 3.1, both $u_{c,h}$ and the operators M, M_h are well defined if $\sigma \geq \sigma_0$, which we assume for the rest of the paper.

Lemma 3.2. *Suppose that $u \in H^s(\Omega)$ and let $u_{c,h} \in V_h$ be defined by (3.3). Then there holds*

$$(3.12) \quad \|u - u_{c,h}\|_{1,h} + h\|u - u_{c,h}\|_{2,h} \lesssim (1 + \sigma)h^{\ell-1} \|u\|_{H^\ell(\Omega)},$$

where $\ell = \min\{k + 1, s\}$ and k denotes the polynomial degree of the finite space V_h .

Proof. By (3.10), (3.3), and (3.9), we have for any $v \in V_h$,

$$\begin{aligned} \|u - u_{c,h}\|_{1,h} &\leq \|u - v\|_{1,h} + \|L_h^{-1}L_h(v - u_{c,h})\|_{1,h} \\ &\lesssim \|u - v\|_{1,h} + \|L(v - u)\|_{-1,h} \\ &\lesssim (1 + \sigma)\|u - v\|_{1,h}. \end{aligned}$$

As v was arbitrary, the first estimate in (3.12) follows from (3.8).

To obtain the second error estimate in (3.12), we use the inverse inequality (3.7) to conclude that

$$\begin{aligned} \|u - u_{c,h}\|_{2,h} &\lesssim \|u - v\|_{2,h} + h^{-1}\|u_{c,h} - v\|_{1,h} \\ &\lesssim \|u - v\|_{2,h} + h^{-1}\|u - v\|_{1,h} + h^{-1}\|u - u_{c,h}\|_{1,h}. \end{aligned}$$

Again, choosing v so that (3.8) holds, we obtain (3.12). □

With the preliminary analysis completed, we turn our attention to the nonlinear problem (2.7) and the mapping M defined in (3.1). Note that by (3.1), (2.3), and (3.3), for all $w \in H^3(\Omega; \mathcal{T}_h)$,

$$\begin{aligned}
 (3.13) \quad Mw &= L_h^{-1}(Lw - Fw) \\
 &= L_h^{-1}(Lw - L(w - u) - R(w - u)) \\
 &= L_h^{-1}(Lu - R(w - u)) \\
 &= u_{c,h} - L_h^{-1}R(w - u),
 \end{aligned}$$

and hence

$$(3.14) \quad Mw_1 - Mw_2 = L_h^{-1}(R(w_2 - u) - R(w_1 - u)) \quad \forall w_1, w_2 \in H^3(\Omega; \mathcal{T}_h).$$

From (3.14) and (3.13) it is clear that the crucial ingredient for applying the Banach fixed point theorem to the map M_h in a small ball around $u_{c,h}$ is the contraction estimate of R established in the next lemma.

Lemma 3.3 (Contraction Estimate of R). *For any $w_1, w_2 \in H^3(\Omega; \mathcal{T}_h)$, there holds*

$$(3.15) \quad \|Rw_1 - Rw_2\|_{-1,h} \lesssim (1 + |\ln h|^{\frac{1}{2}})(\|w_1\|_{2,h} + \|w_2\|_{2,h})\|w_1 - w_2\|_{2,h}.$$

Proof. By (2.5) and (2.8) we have for any $v \in V_h$,

$$\begin{aligned}
 \langle Rw_1 - Rw_2, v \rangle &= \int_{\Omega} (\det(D_h^2 w_2) - \det(D_h^2 w_1)) v \, dx \\
 &+ \sum_{e \in \mathcal{E}_h^i} \int_e \left([[\{\{\text{cof}(D_h^2 w_1)\}\} \nabla w_1]] - [[\{\{\text{cof}(D_h^2 w_2)\}\} \nabla w_2]] \right) v \, ds \\
 &- \sum_{e \in \mathcal{E}_h^b} \int_e \left([[\text{cof}(D_h^2 w_1) \nabla v]] w_1 - [[\text{cof}(D_h^2 w_2) \nabla v]] w_2 \right) ds \\
 &= \frac{1}{2} \int_{\Omega} (\text{cof}(D_h^2(w_1 + w_2)) : D_h^2(w_2 - w_1)) v \, dx \\
 &+ \sum_{e \in \mathcal{E}_h^i} \int_e \left([[\{\{\text{cof}(D_h^2(w_1 - w_2))\}\} \nabla w_1]] - [[\{\{\text{cof}(D_h^2 w_2)\}\} \nabla(w_2 - w_1)]] \right) v \, ds \\
 &- \sum_{e \in \mathcal{E}_h^b} \int_e \left([[\text{cof}(D_h^2 w_1) \nabla v]](w_1 - w_2) - [[\text{cof}(D_h^2(w_2 - w_1)) \nabla v]] w_2 \right) ds.
 \end{aligned}$$

Using the inverse inequality (3.7), (2.9), (3.5), and the Cauchy-Schwarz inequality we obtain

$$\begin{aligned}
 \langle Rw_1 - Rw_2, v \rangle &\leq \left(\frac{1}{2} \sum_{T \in \mathcal{T}_h} |w_1 + w_2|_{H^2(T)} |w_1 - w_2|_{H^2(T)} \right. \\
 &+ \sum_{e \in \mathcal{E}_h^i} \left(\| \{\{D_h^2(w_1 - w_2)\}\} \|_{L^2(e)} \| [[\nabla w_1]] \|_{L^2(e)} \right. \\
 &\quad \left. \left. + \| \{\{D_h^2 w_2\}\} \|_{L^2(e)} \| [[\nabla(w_1 - w_2)]] \|_{L^2(e)} \right) \right) \|v\|_{L^\infty(\Omega)} \\
 &+ \sum_{e \in \mathcal{E}_h^b} \left(\|D_h^2 w_1\|_{L^2(e)} \|w_1 - w_2\|_{L^2(e)} \right)
 \end{aligned}$$

$$\begin{aligned}
 & + \|D_h^2(w_1 - w_2)\|_{L^2(e)} \|w_2\|_{L^2(e)} \|\nabla v\|_{L^\infty(\Omega)} \\
 \lesssim & (1 + |\ln h|^{\frac{1}{2}}) \left((\|w_1\|_{2,h} + \|w_2\|_{2,h}) \|w_1 - w_2\|_{2,h} \right. \\
 & + \sum_{e \in \mathcal{E}_h^i} \left(h_e^{\frac{1}{2}} \|\{D_h^2(w_1 - w_2)\}\|_{L^2(e)} h^{-\frac{1}{2}} \|\llbracket \nabla w_1 \rrbracket\|_{L^2(e)} \right. \\
 & \quad \left. + h^{\frac{1}{2}} \|\{D_h^2 w_2\}\|_{L^2(e)} h^{-\frac{1}{2}} \|\llbracket \nabla(w_1 - w_2) \rrbracket\|_{L^2(e)} \right) \\
 & \quad \left. + \sum_{e \in \mathcal{E}_h^b} h_e^{-1} \left(h^{\frac{1}{2}} \|D_h^2 w_1\|_{L^2(e)} h^{-\frac{1}{2}} \|w_1 - w_2\|_{L^2(e)} \right. \right. \\
 & \quad \left. \left. + h^{\frac{1}{2}} \|D_h^2(w_1 - w_2)\|_{L^2(e)} h^{-\frac{1}{2}} \|w_2\|_{L^2(e)} \right) \right) \|v\|_{1,h} \\
 \lesssim & (1 + |\ln h|^{\frac{1}{2}}) (\|w_1\|_{2,h} + \|w_2\|_{2,h}) \|w_1 - w_2\|_{2,h} \|v\|_{1,h}.
 \end{aligned}$$

The estimate (3.15) then follows from (3.6). □

In light of (3.13)–(3.14), Lemma 3.3 immediately gives us the following three results.

Lemma 3.4 (Contraction Property of M on $H^3(\Omega; \mathcal{T}_h)$). *For any $w_1, w_2 \in H^3(\Omega; \mathcal{T}_h)$ there holds*

$$\begin{aligned}
 (3.16) \quad & \|Mw_1 - Mw_2\|_{1,h} \\
 & \lesssim (1 + |\ln h|^{\frac{1}{2}}) (\|u - w_1\|_{2,h} + \|u - w_2\|_{2,h}) \|w_1 - w_2\|_{2,h}.
 \end{aligned}$$

Proof. By (3.14), (3.10), and (3.15), we have

$$\begin{aligned}
 \|Mw_1 - Mw_2\|_{1,h} & = \|L_h^{-1}(R(w_2 - u) - R(w_1 - u))\|_{1,h} \\
 & \lesssim \|R(w_2 - u) - R(w_1 - u)\|_{-1,h} \\
 & \lesssim (1 + |\ln h|^{\frac{1}{2}}) (\|u - w_1\|_{2,h} + \|u - w_2\|_{2,h}) \|w_1 - w_2\|_{2,h}. \quad \square
 \end{aligned}$$

Lemma 3.5 (Contraction Property of M_h on V_h). *Define the discrete closed ball with center $u_{c,h}$ and radius ρ as*

$$(3.17) \quad \mathbb{B}_\rho(u_{c,h}) = \{v \in V_h; \|u_{c,h} - v\|_{1,h} \leq \rho\}.$$

Then there exists a constant $C_1 > 0$ such that for any $v_1, v_2 \in \mathbb{B}_\rho(u_{c,h})$ there holds

$$\begin{aligned}
 (3.18) \quad & \|M_h v_1 - M_h v_2\|_{1,h} \\
 & \leq C_1 h^{-2} (1 + |\ln h|^{\frac{1}{2}}) (\rho + (1 + \sigma) h^{\ell-1} \|u\|_{H^\ell(\Omega)}) \|v_1 - v_2\|_{1,h}.
 \end{aligned}$$

Proof. By (3.16), (3.17), the inverse inequality (3.7) and (3.12), we have

$$\begin{aligned}
 \|M_h v_1 - M_h v_2\|_{1,h} & \lesssim (1 + |\ln h|^{\frac{1}{2}}) (\|u - v_1\|_{2,h} + \|u - v_2\|_{2,h}) \|v_1 - v_2\|_{2,h} \\
 & \lesssim h^{-1} (1 + |\ln h|^{\frac{1}{2}}) (\|u - u_{c,h}\|_{2,h} + h^{-1} \rho) \|v_1 - v_2\|_{1,h} \\
 & \lesssim h^{-2} (1 + |\ln h|^{\frac{1}{2}}) (\rho + (1 + \sigma) h^{\ell-1} \|u\|_{H^\ell(\Omega)}) \|v_1 - v_2\|_{1,h}. \quad \square
 \end{aligned}$$

Lemma 3.6 (Mapping Property of M_h). *There holds for any $v \in \mathbb{B}_\rho(u_{c,h})$,*

$$(3.19) \quad \|u_{c,h} - M_h v\|_{1,h} \leq C_2 h^{-2} (1 + |\ln h|^{\frac{1}{2}}) (\rho^2 + (1 + \sigma)^2 h^{2(\ell-1)} \|u\|_{H^\ell(\Omega)}^2).$$

Proof. For any $w \in H^3(\Omega; \mathcal{T}_h)$, we have by (3.13), (3.10), and (3.15) that

$$\begin{aligned} \|u_{c,h} - M w\|_{1,h} &= \|L_h^{-1} R(w - u)\|_{1,h} \\ &\lesssim \|R(w - u)\|_{-1,h} \\ &\lesssim (1 + |\ln h|^{\frac{1}{2}}) \|u - w\|_{2,h}^2. \end{aligned}$$

Therefore by the inverse inequality (3.7), (3.17) and (3.12), we have for $v \in \mathbb{B}_\rho(u_{c,h})$,

$$\begin{aligned} \|u_{c,h} - M_h v\|_{1,h} &\lesssim (1 + |\ln h|^{\frac{1}{2}}) \|u - v\|_{2,h}^2 \\ &\lesssim (1 + |\ln h|^{\frac{1}{2}}) (\|u - u_{c,h}\|_{2,h}^2 + \|u_{c,h} - v\|_{2,h}^2) \\ &\lesssim h^{-2} (1 + |\ln h|^{\frac{1}{2}}) ((1 + \sigma)^2 h^{2(\ell-1)} \|u\|_{H^\ell(\Omega)}^2 + \rho^2). \quad \square \end{aligned}$$

With the preceding results established, we are now in position to prove the first main result.

Theorem 3.1 (Main Theorem I). *There exists an $h_0(\sigma) > 0$ such that for $h \leq h_0(\sigma)$ there exists a solution u_h to the penalty method (2.7). Moreover,*

$$(3.20) \quad \|u - u_h\|_{1,h} \lesssim (1 + \sigma) h^{\ell-1} \|u\|_{H^\ell(\Omega)},$$

$$(3.21) \quad \|u - u_h\|_{2,h} \lesssim (1 + \sigma) h^{\ell-2} \|u\|_{H^\ell(\Omega)}.$$

Proof. Since $\ell = \min\{k + 1, s\} > 3$, we can choose $h_0(\sigma) > 0$ such that $h \leq h_0(\sigma)$ implies

$$(3.22) \quad \delta = 2 \max\{C_1, C_2\} h^{-2} (1 + |\ln h|^{\frac{1}{2}}) (1 + \sigma) h^{\ell-1} \|u\|_{H^\ell(\Omega)} < 1.$$

Fix $h \leq h_0(\sigma)$ and set

$$(3.23) \quad \rho_0 = (1 + \sigma) h^{\ell-1} \|u\|_{H^\ell(\Omega)}.$$

Then for any $v \in \mathbb{B}_{\rho_0}(u_{c,h})$, we have by (3.19), (3.22) and (3.23) that

$$\begin{aligned} \|u_{c,h} - M_h v\|_{1,h} &\leq C_2 h^{-2} (1 + |\ln h|^{\frac{1}{2}}) (\rho_0^2 + (1 + \sigma)^2 h^{2(\ell-1)} \|u\|_{H^\ell(\Omega)}^2) \\ &= 2 \left(C_2 h^{-2} (1 + |\ln h|^{\frac{1}{2}}) (1 + \sigma) h^{\ell-1} \|u\|_{H^\ell(\Omega)} \right) \rho_0 \\ &\leq \rho_0, \end{aligned}$$

and so M_h maps $\mathbb{B}_{\rho_0}(u_{c,h})$ into $\mathbb{B}_{\rho_0}(u_{c,h})$. Moreover, by (3.18) and (3.22) for $v_1, v_2 \in \mathbb{B}_{\rho_0}(u_{c,h})$,

$$\begin{aligned} \|M_h v_1 - M_h v_2\|_{1,h} &\leq C_1 h^{-1} (1 + |\ln h|^{\frac{1}{2}}) (\rho_0 + (1 + \sigma) h^{\ell-1} \|u\|_{H^\ell(\Omega)}) \|v_1 - v_2\|_{1,h} \\ &= 2 C_1 h^{-1} (1 + |\ln h|^{\frac{1}{2}}) (1 + \sigma) h^{\ell-1} \|u\|_{H^\ell(\Omega)} \|v_1 - v_2\|_{1,h} \\ &\leq \delta \|v_1 - v_2\|_{1,h}. \end{aligned}$$

Hence M_h has a unique fixed point $u_h \in \mathbb{B}_{\rho_0}(u_{c,h})$, which is a solution of (2.7). Moreover, by (3.23) and (3.12),

$$\|u - u_h\|_{1,h} \leq \|u - u_{c,h}\|_{1,h} + \rho_0 \lesssim (1 + \sigma) h^{\ell-1} \|u\|_{H^\ell(\Omega)}.$$

Finally, by the inverse inequality, we have

$$\|u - u_h\|_{2,h} \leq \|u - u_{c,h}\|_{2,h} + h^{-1} \rho_0 \lesssim (1 + \sigma) h^{\ell-2} \|u\|_{H^\ell(\Omega)}. \quad \square$$

Remark 3.4. Theorem 3.1 shows that $\|u_h - u_{c,h}\|_{1,h} \leq \rho_0 = (1 + \sigma)h^{\ell-1}\|u\|_{H^\ell(\Omega)}$, but this estimate is not sharp. Indeed, by (3.13) we have

$$(3.24) \quad u_h = M_h u_h = u_{c,h} - L_h^{-1}R(u_h - u),$$

and therefore by (3.1), (3.15), and (3.21),

$$(3.25) \quad \begin{aligned} \|u_h - u_{c,h}\|_{1,h} &= \|L_h^{-1}R(u_h - u)\|_{1,h} \\ &\lesssim \|R(u_h - u)\|_{-1,h} \\ &\lesssim (1 + |\ln h|^{\frac{1}{2}})\|u - u_h\|_{2,h}^2 \\ &\lesssim (1 + |\ln h|^{\frac{1}{2}})(1 + \sigma)^2 h^{2(\ell-2)}\|u\|_{H^\ell(\Omega)}^2. \end{aligned}$$

Theorem 3.2 (Main Theorem II). *In addition to the hypotheses of Theorem 3.1, assume that $u \in W^{3,\infty}(\Omega)$. Then there holds*

$$(3.26) \quad \|u - u_h\|_{L^2(\Omega)} \lesssim (1 + \sigma)^2 (h^\ell \|u\|_{H^\ell(\Omega)} + (1 + |\ln h|^{\frac{1}{2}})h^{2(\ell-2)}\|u\|_{H^\ell(\Omega)}^2).$$

Proof. Let $\psi \in H_0^1(\Omega)$ be the solution to the following auxiliary problem:

$$(3.27a) \quad -\nabla \cdot (\text{cof}(D^2 u) \nabla \psi) = u - u_h \quad \text{in } \Omega,$$

$$(3.27b) \quad \psi = 0 \quad \text{on } \partial\Omega.$$

Since $u \in W^{3,\infty}(\Omega)$, there holds $\text{cof}(D^2 u) \in [W^{1,\infty}(\Omega)]^{2 \times 2}$. Therefore by elliptic regularity theory [23, 17, 24], we have

$$(3.28) \quad \|\psi\|_{H^2(\Omega)} \lesssim \|u - u_h\|_{L^2(\Omega)}.$$

Let $\psi_h \in V_h$ be chosen so that

$$(3.29) \quad \|\psi - \psi_h\|_{1,h} \lesssim h\|\psi\|_{H^2(\Omega)} \lesssim h\|u - u_h\|_{L^2(\Omega)}.$$

Note that by the trace inequality,

$$\left(\sum_{e \in \mathcal{E}_h^b} h_e \|\nabla \psi\|_{L^2(e)}^2 \right)^{\frac{1}{2}} \lesssim \|\psi\|_{H^2(\Omega)},$$

and therefore by (3.4) and (3.27b),

$$\|\psi\|_{1,h} \lesssim \|\psi\|_{H^2(\Omega)}.$$

Hence, by (3.29) we have

$$(3.30) \quad \|\psi_h\|_{1,h} \leq \|\psi_h - \psi\|_{1,h} + \|\psi\|_{1,h} \lesssim \|\psi\|_{H^2(\Omega)} \lesssim \|u - u_h\|_{L^2(\Omega)}.$$

Using (3.27), we can write

$$(3.31) \quad \begin{aligned} \|u - u_h\|_{L^2(\Omega)}^2 &= \langle L(u - u_h), \psi \rangle \\ &= \langle L(u - u_h), \psi - \psi_h \rangle + \langle L(u - u_h), \psi_h \rangle. \end{aligned}$$

From (3.9), (3.29) and (3.28), we have

$$(3.32) \quad \begin{aligned} \langle L(u - u_h), \psi - \psi_h \rangle &\leq \|L(u - u_h)\|_{-1,h} \|\psi - \psi_h\|_{1,h} \\ &\lesssim (1 + \sigma)\|u - u_h\|_{1,h} \|\psi - \psi_h\|_{1,h} \\ &\lesssim h(1 + \sigma)\|u - u_h\|_{1,h} \|u - u_h\|_{L^2(\Omega)}. \end{aligned}$$

On the other hand, (3.24), (3.15), (3.30) and (3.28) imply

$$\begin{aligned}
 (3.33) \quad \langle L(u_h - u), \psi_h \rangle &= \langle L_h(u_h - u_{c,h}), \psi_h \rangle \\
 &= \langle R(u - u_h), \psi_h \rangle \\
 &\leq \|R(u - u_h)\|_{-1,h} \|\psi_h\|_{1,h} \\
 &\lesssim (1 + |\ln h|^{\frac{1}{2}}) \|u - u_h\|_{2,h}^2 \|\psi_h\|_{1,h} \\
 &\lesssim (1 + |\ln h|^{\frac{1}{2}}) \|u - u_h\|_{2,h}^2 \|u - u_h\|_{L^2(\Omega)}.
 \end{aligned}$$

Combining (3.31)–(3.33), dividing by $\|u - u_h\|_{L^2(\Omega)}$, and using (3.20)–(3.21), we obtain

$$\|u - u_h\|_{L^2(\Omega)} \lesssim (1 + \sigma)^2 (h^\ell \|u\|_{H^\ell(\Omega)} + (1 + |\ln h|^{\frac{1}{2}}) h^{2(\ell-2)} \|u\|_{H^\ell(\Omega)}^2). \quad \square$$

Remark 3.5. Since $\ell > 3$, the error estimate (3.26) is of higher order than (3.20). Moreover, the L^2 estimate is (almost) of the optimal order $k + 1$ provided $s \geq 4$.

4. NUMERICAL EXPERIMENTS

In this section, we perform some numerical tests that back up the theoretical results proved in the previous section, as well as show the effectiveness and efficiency of the method. We solve the finite element method (2.7) using the COMSOL Multiphysics software package [13], and solve the resulting nonlinear algebraic system using Newton’s method. Given $u_0 \in V_h$, the Newton approximations to u_h form a sequence $\{u_k\}_{k=0}^\infty \subset V_h$ satisfying

$$(4.1) \quad DF_h[u_k](u_{k+1} - u_k) = -F_h u_k,$$

where $DF_h : V_h \rightarrow \mathcal{L}(V_h; V'_h)$ denotes the Gâteaux derivative of F_h . Here, $\mathcal{L}(V_h; V'_h)$ denotes the space of linear maps from V_h to V'_h .

Similar to other numerical methods for the Monge-Ampère equation [22, 26], the Newton iteration requires an accurate starting value to ensure convergence. In order to obtain a good initial guess, we apply the vanishing moment methodology to the Monge-Ampère equation [18]. The crux of the vanishing moment method is to approximate fully nonlinear PDEs by higher order semi-linear PDEs, in particular, fourth order PDEs. For the case of the Monge-Ampère equation (1.1) the vanishing moment approximation is defined to be the solution to the following fourth order problem:

$$(4.2) \quad -\varepsilon \Delta^2 u^\varepsilon + \det(D^2 u^\varepsilon) = f \quad 0 < \varepsilon \ll 1,$$

along with appropriate boundary conditions. We refer the reader to [18, 21] for more details and motivation of the vanishing moment method. Besides being a valuable tool to obtain close approximations of the Monge-Ampère equation, another interesting feature of the vanishing moment methodology is the ability to choose the concave solution (recall there is a convex solution and concave solution to the Monge-Ampère equation) by substituting ε by $-\varepsilon$ in (4.2). However, we do not pursue this direction in the numerical experiments below.

The C^0 interior penalty method for (4.2) is defined as seeking $u_h^\varepsilon \in V_h$ such that

$$(4.3) \quad \varepsilon A_h u_h^\varepsilon + F_h u_h^\varepsilon = 0,$$

TABLE 1. Example 1. Errors of computed solution and rates of convergence with respect to h with $\sigma = 100$.

	h	$\ u - u_h\ _{L^2(\Omega)}$	rate	$ u - u_h _{H^1(\Omega)}$	rate	$ u - u_h _{H^2(\mathcal{T}_h)}$	rate
$k = 2$	1/8	3.06E-03		7.69E-02		7.88E+00	
	1/16	1.62E-03	0.92	4.74E-02	0.70	6.33E+00	0.31
	1/32	1.96E-04	3.05	1.13E-02	2.06	3.18E+00	0.99
	1/64	2.79E-05	2.81	2.77E-03	2.03	1.59E+00	1.00
	1/128	7.07E-06	1.98	6.88E-04	2.01	7.97E-01	1.00
	1/256	2.03E-06	1.80	1.72E-04	2.00	3.99E-01	1.00
$k = 3$	1/8	1.47E-04		1.87E-03		4.27E-01	
	1/16	5.62E-05	1.38	9.45E-04	0.98	2.72E-01	0.65
	1/32	3.89E-06	3.85	9.51E-05	3.31	6.82E-02	2.00
	1/64	2.55E-07	3.93	1.00E-05	3.25	1.71E-02	2.00
	1/128	1.64E-08	3.96	1.13E-06	3.14	4.27E-03	2.00
	1/256	1.08E-09	3.92	1.35E-07	3.07	1.07E-03	2.00
$k = 4$	1/8	5.21E-06		7.99E-05		1.23E-02	
	1/16	1.68E-06	1.63	3.09E-05	1.37	6.22E-03	0.99
	1/32	6.06E-08	4.79	1.79E-06	4.11	7.80E-04	3.00
	1/64	2.04E-09	4.89	1.06E-07	4.07	9.72E-05	3.01
	1/128	6.47E-11	4.98	6.52E-09	4.03	1.22E-05	3.00

where

$$\begin{aligned} \langle A_h v, w \rangle = & \int_{\Omega} D_h^2 v : D_h^2 w \, dx - \sum_{e \in \mathcal{E}_h^i} \int_e \left(\{\{\partial_{nn}^2 v\}\} [[\nabla w]] + [[\nabla v]] \{\{\partial_{nn}^2 w\}\} \right. \\ & \left. - \sigma h_e^{-1} [[\nabla v]] [[\nabla w]] \right) ds \quad \forall v, w \in V_h, \end{aligned}$$

and

$$\{\{\partial_{nn}^2 w\}\}|_e = \frac{1}{2} (D_h^2 w^+ \mathbf{n}_+ \cdot \mathbf{n}_+|_e + D_h^2 w^- \mathbf{n}_- \cdot \mathbf{n}_-|_e) \quad e = \partial T^+ \cap \partial T^- \in \mathcal{E}_h^i$$

denotes the average of the second order normal derivative of w . To solve (4.3), we again use Newton’s method, which creates a sequence $\{u_k^\varepsilon\}_{k=0}^\infty \subset V_h$ satisfying

$$(4.4) \quad \varepsilon A_h u_{k+1}^\varepsilon + DF_h[u_k^\varepsilon](u_{k+1}^\varepsilon - u_k^\varepsilon) = -F_h u_k^\varepsilon.$$

We plan to address the convergence of the Newton iterations (4.1) and (4.4) in the near future. For now, we provide some numerical experiments illustrating the robustness of this technique.

4.1. Example 1. In this test, we solve (2.7) for varying values of h and k , and choose our data such that the exact solution to the Monge-Ampère equation (1.1) is $u = 20e^{x_1^6/6+x_2}$. We take $\Omega = (0, 1)^2$, the unit square and set $\sigma = 100$.

In order to obtain some good initial guesses, we first solve (4.3) with ε -values 1E-2, 1E-4, 1E-6, and 0, using each previous solution as our initial guess in the Newton iteration (4.4) (we take $u_0^\varepsilon = x_1^2 + x_2^2$ as our initial guess for the first iteration with $\varepsilon = 1E-2$). After computing the solution of (2.7) we calculate the

L^2 , H^1 , and broken H^2 error and record the errors in Table 1¹. As predicted by the theoretical results in Theorem 3.1, for $k \geq 3$ we observe $(k + 1)$, (k) , and $(k - 1)$ -order convergence in the L^2 , H^1 and H^2 -norms. Furthermore, the numerical tests also indicate that the method is convergent using quadratic polynomials, although a theoretical proof of such a result has yet to be shown.

4.2. Example 2. In this test, we solve the Monge-Ampère equation on the unit disc $\Omega = B(0,1)$ with the same test problem as Example 1. Our finite element space is constructed such that on curved elements, we use polynomial functions of degree $\leq k$ in the curvilinear coordinates for T , which in this case, are isoparametric/subparametric finite elements.

We solve (2.7) for varying h and k values and record the errors in Table 2. In order to obtain good initial guesses for the Newton iteration, we use the same strategy as in Example 1, first solving the regularized problem (4.3). As expected, the convergence rates are exactly the same as the previous test with convergence rates of $O(h^4)$, $O(h^3)$, $O(h^2)$ in the L^2 , H^1 , and H^2 -norms using cubic polynomials. Similar to the previous tests, we also observe that the method is also convergent using quadratic polynomials.

TABLE 2. Example 2. Errors of computed solution and rates of convergence with respect to h with $\sigma = 100$ on curved domain.

	h	$\ u - u_h\ _{L^2(\Omega)}$	rate	$ u - u_h _{H^1(\Omega)}$	rate	$ u - u_h _{H^2(\mathcal{T}_h)}$	rate
$k = 2$	1/8	3.30E-03		1.35E-01		9.73E+00	
	1/16	1.03E-03	1.68	3.50E-02	1.95	4.90E+00	0.99
	1/32	1.66E-04	2.64	1.01E-02	1.79	2.41E+00	1.02
	1/64	3.06E-05	2.44	1.89E-03	2.42	1.12E+00	1.10
	1/128	1.11E-05	1.46	4.66E-04	2.02	5.54E-01	1.02
$k = 3$	1/8	9.40E-04		4.91E-03		8.73E-01	
	1/16	7.85E-05	3.58	5.82E-04	3.08	2.24E-01	1.96
	1/32	4.63E-06	4.08	5.66E-05	3.36	4.89E-02	2.20
	1/64	2.99E-07	3.95	6.51E-06	3.12	1.16E-02	2.07
	1/128	1.89E-08	3.98	7.91E-07	3.04	2.84E-03	2.04

4.3. Example 3. For the last test we solve (2.7) on the unit square $\Omega = (0, 1)^2$ using quadratic polynomials and choose our data such that the exact solution is $u = \frac{(4(x_1^2 + x_2^2))^{\frac{3}{4}}}{3}$. Unlike the previous two tests, the exact solution is not smooth, and one can readily check that $u \in W^{2,p}(\Omega)$ for all $p \in [1, 4)$ but $u \notin C^2(\bar{\Omega})$ [15].

Using the same strategy as the previous two tests, we first solve the regularized problem (4.3) with $\sigma = 100$ using quadratic polynomials to obtain some good initial guesses. We then solve (2.7) and record the errors in Table 3. Although a theoretical proof of such a result has not been proven, the numerical experiments indicate that the method is convergent with rates of $O(h^2)$, $O(h^{\frac{3}{2}})$, and $O(h^{\frac{1}{2}})$ in the L^2 , H^1 , and H^2 -norms. The suboptimal rates are most likely due to the low regularity of the solution.

¹We define the piecewise H^2 -seminorm as $|u - u_h|_{H^2(\mathcal{T}_h)} = \left(\sum_{T \in \mathcal{T}_h} |u - u_h|_{H^2(T)}^2\right)^{\frac{1}{2}}$.

TABLE 3. Example 3. Errors of computed solution and rates of convergence with respect to h with $\sigma = 100$ for nonsmooth solution.

	h	$\ u - u_h\ _{L^2(\Omega)}$	rate	$ u - u_h _{H^1(\Omega)}$	rate	$ u - u_h _{H^2(\mathcal{T}_h)}$	rate
$k = 2$	1/8	2.96E-05		2.61E-03		2.52E-01	
	1/16	6.33E-06	2.22	9.46E-04	1.46	1.79E-01	0.49
	1/32	1.48E-06	2.09	3.44E-04	1.46	1.27E-01	0.49
	1/64	3.73E-07	1.99	1.26E-04	1.45	9.03E-02	0.50
	1/128	9.74E-08	1.94	4.66E-05	1.43	6.41E-02	0.50

5. CONCLUDING REMARKS

In this paper, we have developed and analyzed C^0 penalty methods for the fully nonlinear Monge-Ampère equation. To build convergent numerical schemes, we constructed discretizations such that the resulting discrete linearization is stable, symmetric, and consistent with the continuous linearization. With this in hand, we proved existence of the numerical solution as well as derived quasi-optimal error estimates using a simple fixed-point technique.

The methodology developed in this paper has been applied to the three dimensional Monge-Ampère equation [8]. Due to the simplicity and flexibility of the method, it is also relatively straightforward to formulate discretizations for more general Monge-Ampère equations in which the function f depends on ∇u and u . We plan to address the convergence properties of such problems in the near future.

REFERENCES

- [1] F. E. Baginski and N. Whitaker, *Numerical solutions of boundary value problems for K -surfaces in R^3* , Numer. Methods for PDEs, 12(4):525–546, 1996. MR1396470 (97h:65152)
- [2] G. Barles and P. E. Souganidis, *Convergence of approximation schemes for fully nonlinear second order equations*, Asymptotic Anal., 4(3):271–283, 1991. MR1115933 (92d:35137)
- [3] J. D. Benamou, B. D. Froese, and A. M. Oberman, *Two numerical methods for the elliptic Monge-Ampère equation*, M2AN Math. Model. Numer. Anal., DOI 10.1051/m2an/2010017, 2010.
- [4] C. Bernardi, *Optimal finite element interpolation on curved domains*, SIAM J. Numer. Anal., 26(5):1212–1240, 1989. MR1014883 (91a:65228)
- [5] K. Böhmer, *On finite element methods for fully nonlinear elliptic equations of second order*, SIAM J. Numer. Anal., 46(3):1212–1249, 2008. MR2390991 (2009b:35092)
- [6] J. H. Bramble, J. E. Pasciak, and A. H. Schatz, *The construction of preconditioners for elliptic problems by substructuring. I.*, Math. Comp. 47(175):103–134, 1986. MR842125 (87m:65174)
- [7] S. C. Brenner and L. R. Scott, *The Mathematical Theory of Finite Element Methods*, third edition, Springer, 2008. MR2373954 (2008m:65001)
- [8] S. C. Brenner and M. Neilan, *Finite element approximations of the three dimensional Monge-Ampère equation*, submitted.
- [9] C. J. Budd, W. Huang, and R. D. Russell, *Adaptivity with moving grids*, Acta Numerica 18:111–241, 2009. MR2506041 (2010c:65165)
- [10] L. A. Caffarelli, L. Nirenberg, and J. Spruck, *The Dirichlet problem for nonlinear second-order elliptic equations I. Monge-Ampère equation*, Comm. Pure Appl. Math. 37(3):369–402, 1984. MR739925 (87f:35096)
- [11] L. A. Caffarelli and M. Milman, *Monge-Ampère Equation: Applications to Geometry and Optimization*, American Mathematical Society, Providence, RI, 1999. MR1660738 (99f:00018)
- [12] P. G. Ciarlet, *The Finite Element Method for Elliptic Problems*. North-Holland, Amsterdam, 1978. MR0520174 (58:25001)
- [13] COMSOL Multiphysics Software Package, <http://www.comsol.com>.

- [14] R. Courant and D. Hilbert, *Methods of Mathematical Physics*, John Wiley & Sons, Inc., New York, 1989. MR1013360 (90k:35001)
- [15] E. J. Dean and R. Glowinski, *Numerical methods for fully nonlinear elliptic equations of the Monge-Ampère type*, *Comput. Methods Appl. Mech. Engrg.*, 195(13-16):1344–1386, 2006. MR2203972 (2006i:65191)
- [16] G. L. Delzanno, L. Chacón, J. M. Finn, Y. Chung, and G. Lapenta, *An optimal robust equidistribution method for two-dimensional grid adaptation based on Monge-Kantorovich optimization*, *J. Comput. Phys.* 227(23):9841–9864, 2008. MR2469037 (2010b:65286)
- [17] L. C. Evans, *Partial Differential Equations*, volume 19 of *Graduate Studies in Mathematics*, American Mathematical Society, Providence, RI, 1998. MR1625845 (99e:35001)
- [18] X. Feng and M. Neilan, *Vanishing moment method and moment solutions for second order fully nonlinear partial differential equations*, *J. Scient. Comp.*, 38(1):74–98, 2009. MR2472219 (2010a:65234)
- [19] X. Feng and M. Neilan, *Mixed finite element methods for the fully nonlinear Monge-Ampère equation based on the vanishing moment method*, *SIAM J. Numer. Anal.* 47(2):1226–1250, 2009. MR2485451 (2010a:65235)
- [20] X. Feng and M. Neilan, *Analysis of Galerkin methods for the fully nonlinear Monge-Ampère equation*, *J. Sci. Comput.*, DOI: 10.107/S10915-010-9439-1.
- [21] X. Feng and M. Neilan, *The vanishing moment method for fully nonlinear second order partial differential equations: formulation, theory, and numerical analysis*, submitted.
- [22] B. D. Froese and A. M. Oberman, *Convergent finite difference solvers for viscosity solutions of the elliptic Monge-Ampère equation in dimensions two and higher*, arXiv:1007.0765, 2010.
- [23] D. Gilbarg and N. S. Trudinger, *Elliptic Partial Differential Equations of Second Order*, Springer-Verlag, Berlin, 2001. MR1814364 (2001k:35004)
- [24] P. Grisvard, *Elliptic Problems on Nonsmooth Domains*, Pitman Publishing Inc., 1985. MR775683 (86m:35044)
- [25] C. E. Gutiérrez, *The Monge-Ampère Equation*, volume 44 of *Progress in Nonlinear Differential Equations and Their Applications*, Birkhauser, Boston, MA, 2001. MR1829162 (2002e:35075)
- [26] G. Loeper and F. Rapetti, *Numerical solution of the Monge-Ampère equation by a Newton’s algorithm*, *C. R. Math. Acad. Sci. Paris* 340(4):319–324, 2005. MR2121899
- [27] M. Neilan, *A nonconforming Morley finite element method for the fully nonlinear Monge-Ampère equation*, *Numer. Math.*, 115(3):371–394, 2010.
- [28] J. A. Nitsche, *Über ein Variationsprinzip zur Lösung Dirichlet-Problemen bei Verwendung von Teilräumen, die keinen Randbedingungen unterworfen sind*, *Abh. Math. Sem. Univ. Hamburg*, 36:9–15, 1971. MR0341903 (49:6649)
- [29] A. M. Oberman, *Wide stencil finite difference schemes for the elliptic Monge-Ampère equation and functions of the eigenvalues of the Hessian*, *Discrete Contin. Dyn. Syst. Ser. B*, 10(1):221–238, 2008. MR2399429 (2009f:35101)
- [30] V. I. Oliker and L. D. Prussner, *On the numerical solution of the equation $\frac{\partial^2 z}{\partial x^2} \frac{\partial^2 z}{\partial y^2} - (\frac{\partial^2 z}{\partial x \partial y})^2 = f$ and its discretizations*, *Numer. Math.*, 54:271–293, 1988. MR971703 (90h:65164)
- [31] C. Villani, *Topics in Optimal Transportation*, volume 58 of *Graduate Studies in Mathematics*, American Mathematical Society, Providence, RI, 2003. MR1964483 (2004e:90003)

DEPARTMENT OF MATHEMATICS AND CENTER FOR COMPUTATION AND TECHNOLOGY, LOUISIANA STATE UNIVERSITY, BATON ROUGE, LOUISIANA 70803
E-mail address: brenner@math.lsu.edu

DEPARTMENT OF MATHEMATICS, INDIAN INSTITUTE OF SCIENCE, BANGALORE, 560012
E-mail address: gudi@math.iisc.ernet.in

DEPARTMENT OF MATHEMATICS AND CENTER FOR COMPUTATION AND TECHNOLOGY, LOUISIANA STATE UNIVERSITY, BATON ROUGE, LOUISIANA 70803
E-mail address: neilan@math.lsu.edu

DEPARTMENT OF MATHEMATICS AND CENTER FOR COMPUTATION AND TECHNOLOGY, LOUISIANA STATE UNIVERSITY, BATON ROUGE, LOUISIANA 70803
E-mail address: sung@math.lsu.edu