

HYBRIDIZABLE DISCONTINUOUS GALERKIN METHOD FOR ELLIPTIC EQUATIONS WITH NONLINEAR COEFFICIENTS

MINAM MOON¹

¹DEPARTMENT OF MATHEMATICS, KOREA MILITARY ACADEMY, SOUTH KOREA
Email address: hereandnow@kma.ac.kr

ABSTRACT. In this paper, we analyze the hybridizable discontinuous Galerkin (HDG) method for second-order elliptic equations with nonlinear coefficients, which are used in many fields. We present the HDG method that uses a mixed formulation based on numerical trace and flux. Under assumptions on the nonlinear coefficient and H^2 -regularity for a dual problem, we prove that the discrete systems are well-posed and the numerical solutions have the optimal order of convergence as a mesh parameter. Also, we provide a matrix formulation that can be calculated using an iterative technique for numerical experiments. Finally, we present representative numerical examples in 2D to verify the validity of the proof of Theorem 3.10.

1. INTRODUCTION

Most of the practical problems we are interested in are represented by nonlinear PDEs. In particular, problems such as flows in porous media can be modeled with nonlinear PDEs in [1]. A large amount of research in numerous aspects of simulation of the nonlinear PDEs has been in the spotlight to figure out the complexity of nonlinearity and reduce expensive computational costs. Among them, simulation and analysis using the finite element method (FEM) are being actively conducted. In this paper, we provide the analysis of a robust approximation solution on the hybridizable discontinuous Galerkin (HDG) method for the following second-order elliptic equation with a nonlinear coefficient :

$$-\nabla \cdot (\kappa(x, u) \nabla u(x)) = f(x) \quad \text{in } \Omega, \quad (1.1a)$$

$$u(x) = 0 \quad \text{on } \partial\Omega, \quad (1.1b)$$

where $\kappa(x, u)$ is the nonlinear term, $f \in L^2(\Omega)$, and Ω is a bounded polyhedron in \mathbb{R}^n , $n = 2, 3$.

Since Cockburn et al. formally developed the HDG method for second-order elliptic problems in [2], it has been used to give efficient and robust approximate solutions. The HDG

Received by the editors July 12 2022; Revised November 18 2022; Accepted in revised form December 12 2022; Published online December 25 2022.

2010 *Mathematics Subject Classification.* 65N12, 65N15, 65N30.

Key words and phrases. Nonlinear problems, Hybridizable discontinuous Galerkin (HDG), elliptic equation, error analysis.

method retains the advantages of the discontinuous Galerkin (DG) methods such as flexibility in meshing and preserving local conservation of physical quantities and overcoming to shortcomings of the DG by reducing the globally coupled degree of freedom in [3]. Also, the HDG method can outperform continuous Galerkin (CG) methods in various aspects such as parallel computation computing time, floating-point operation count and superconvergence in [4, 5, 6, 7, 8].

Due to these characteristics, the HDG method has been applied and studied for a variety of problems. The optimal order of convergence for elliptic problems was mathematically proven by introducing HDG projection in [9, 10]. Also, based on superconvergence of the HDG method, local postprocessing was developed for linear convection-diffusion equations in [11], for nonlinear cases in [12]. The HDG method was also applied to more substantive problems such as parabolic equations in [13], acoustic wave equations in [14, 15], Stokes equations in [16, 17], compressible Navier-Stokes equations in [18], and incompressible Navier-Stokes equations in [19, 20, 21].

Mathematical analysis of nonlinear problems has been studied with various FEMs. Convergent or superconvergent results for nonlinear elliptic problems were derived based on mixed method in [22, 23], DG method in [24], and multiscale mortar method in [25]. Although it is essential to understand how accurate the approximate solution is, the mathematical analysis of nonlinear problems is limited. This is because, in contrast to linear problems, nonlinear problems are not guaranteed to have well-posedness and are difficult to deal with nonlinear terms for analysis due to its complexity. So, we conduct error estimations for the nonlinear problem (1.1) to investigate how accurate our approximation in the HDG method is.

In this paper, we analyze the HDG method for elliptic equations with nonlinear coefficients and present numerical results. To deal with nonlinear problems, we first assume that a nonlinear coefficient satisfies the Lipschitz continuous and boundedness. We also assume that domains permit the H^2 -regularity estimates to get high order of convergence in [10]. Finally, we assume that the local stabilization parameter is positive. The well-posedness and accuracy depend on the stabilization parameter. Under these assumptions, we will show the well-posedness of the HDG system and give error estimates. We first derive error equations based on the HDG projection, and then get several identities and an upper bound for the projection errors by using error equations. We can derive the optimal convergence ratio of the HDG method for elliptic problems with nonlinear coefficients. Also, to solve the nonlinear problem numerically, we provide a matrix formulation which modifies the technique introduced in [26, 27]. The matrix formulation can be solved using the iterative technique. In numerical experiment, we show error plots for various nonlinear coefficients and check whether the optimal order of convergence is consistent with the mathematical proof.

The paper is organized in the following way: after this introduction, we will define notations regarding functional space, finite element space, and the HDG projection that will be used throughout the paper, and along with them, we give the HDG method for the nonlinear problem. Then in Section 3, we will provide assumptions used in analysis and prove the well-posedness and error estimations for the HDG method. In Section 4, we present the matrix formulation

for implementation and the numerical results to verify the optimal order of convergence. A conclusion is given in Section 5.

2. PRELIMINARIES

In this section, we present the HDG method for solving second order elliptic equations with nonlinear coefficients. For this, we will first introduce some notations of functional spaces, finite element spaces, and HDG projections. These contents and notations will be borrowed from [10].

Let \mathcal{T}_h be a conforming, shape-regular simplicial triangulation of Ω with maximum element diameter of h . We call F an interface of the triangulation \mathcal{T}_h if F either is shared by two neighboring triangles, T_1 and T_2 in \mathcal{T}_h ($F = \overline{T_1} \cap \overline{T_2}$), or is on the boundary $\partial\Omega$ ($F = \overline{T} \cap \partial\Omega$). Let \mathcal{E}_h denote the set of all interfaces of the triangulation \mathcal{T}_h . Note that any interface F lies on the boundary of some triangle ∂T . Set $\partial\mathcal{T}_h = \cup_{T \in \mathcal{T}_h} \partial T$.

We will use the standard notations for Sobolev spaces and their norms on the domain Ω and their boundaries. For example, $\|v\|_{s,\Omega}$, $|v|_{s,\Omega}$, $\|v\|_{s,\partial\Omega}$, $|v|_{s,\partial\Omega}$, $s > 0$, denote the Sobolev norms and semi-norms on Ω and its boundary $\partial\Omega$. For an integer s , the Sobolev spaces H^s are Hilbert spaces and the norms are defined by the L^2 -norms of their weak derivatives up to order s . For a non-integer s , the spaces are defined by interpolation in [28]. When $s = 0$ we will use $\|v\|_\Omega$ instead of $\|v\|_{0,\Omega}$. $\|v\|_\infty$ will denote the standard L^∞ -norm. $\mathbf{H}_{div}(\Omega)$ denote the space of vector-functions with components in $L^2(\Omega)$ and weak divergence in $L^2(\Omega)$. Also, $H^s(\mathcal{T}_h)$, $\mathbf{H}^s(\mathcal{T}_h)$, and $\mathbf{H}_{div}(\mathcal{T}_h)$ are defined as follows:

$$H^s(\mathcal{T}_h) := \prod_{T \in \mathcal{T}_h} H^s(T), \quad \mathbf{H}^s(\mathcal{T}_h) := \prod_{T \in \mathcal{T}_h} (H^s(T))^n, \quad \mathbf{H}_{div}(\mathcal{T}_h) := \prod_{T \in \mathcal{T}_h} \mathbf{H}_{div}(T).$$

For any element $T \in \mathcal{T}_h$ and any interface $F \in \mathcal{E}_h$, we define

$$W(T) := \mathcal{P}_k(T), \quad \mathbf{V}(T) := \mathcal{P}_k(T), \quad M(F) := \mathcal{P}_k(F),$$

where $\mathcal{P}_k(D)$ denotes the set of polynomials of degree at most k on a domain D and $\mathcal{P}_k(T) = (\mathcal{P}_k(T))^n$. Now, we consider the following finite element spaces:

$$\begin{aligned} W_h &:= \{w \in L^2(\mathcal{T}_h) : w|_T \in W(T) \text{ for all } T \in \mathcal{T}_h\}, \\ \mathbf{V}_h &:= \{\mathbf{v} \in \mathbf{L}^2(\mathcal{T}_h) : \mathbf{v}|_T \in \mathbf{V}(T) \text{ for all } T \in \mathcal{T}_h\}, \\ M_h &:= \{\mu \in L^2(\mathcal{E}_h) : \mu|_F \in M(F) \text{ for all } F \in \mathcal{E}_h\}, \end{aligned}$$

where $L^2(\mathcal{T}_h) = \prod_{T \in \mathcal{T}_h} L^2(T)$, $\mathbf{L}^2(\mathcal{T}_h) = (L^2(\mathcal{T}_h))^n$, and $L^2(\mathcal{E}_h) = \prod_{F \in \mathcal{E}_h} L^2(F)$.

For a domain D is a subset of \mathbb{R}^n and scalar-valued functions $u, v \in L^2(D)$, let $(u, v)_D = \int_D uv \, dx$. For vector-valued functions $\mathbf{u}, \mathbf{v} \in \mathbf{L}^2(D)$, we define $(\mathbf{u}, \mathbf{v})_D = \int_D \mathbf{u} \cdot \mathbf{v} \, dx$. For the boundary $\partial D \subset \mathbb{R}^{n-1}$ of D , we define $\langle u, v \rangle_{\partial D} = \int_{\partial D} uv \, ds$. Then, we introduce the following notation:

$$(w, v)_{\mathcal{T}_h} = \sum_{T \in \mathcal{T}_h} (w, v)_T, \quad \langle w, v \rangle_{\partial\mathcal{T}_h} = \sum_{T \in \mathcal{T}_h} \langle w, v \rangle_{\partial T}.$$

Since the HDG method was developed based on the mixed method and the DG method, We need to consider the following mixed formulation:

$$\alpha(u)\mathbf{q} + \nabla u = 0 \quad \text{in } \Omega, \quad (2.1a)$$

$$\nabla \cdot \mathbf{q} = f \quad \text{in } \Omega, \quad (2.1b)$$

$$u = 0 \quad \text{on } \partial\Omega, \quad (2.1c)$$

where $\alpha(u) = \kappa(u)^{-1}$.

Based on the mixed formulation (2.1), we need to find an approximation to $(u, \mathbf{q}, u|_{\varepsilon_h})$ by the HDG method. For this, the HDG method provides approximations $(u_h, \mathbf{q}_h, \hat{u}_h) \in W_h \times \mathbf{V}_h \times M_h$, determined by the following five equations:

For any $(w, \mathbf{v}, \mu) \in W_h \times \mathbf{V}_h \times M_h$, we require

$$(\alpha(u_h)\mathbf{q}_h, \mathbf{v})_{\mathcal{T}_h} - (u_h, \nabla \cdot \mathbf{v})_{\mathcal{T}_h} + \langle \hat{u}_h, \mathbf{v} \cdot \mathbf{n} \rangle_{\partial\mathcal{T}_h} = 0, \quad (2.2a)$$

$$-(\mathbf{q}_h, \nabla w)_{\mathcal{T}_h} + \langle \hat{\mathbf{q}}_h \cdot \mathbf{n}, w \rangle_{\partial\mathcal{T}_h} = (f, w)_{\mathcal{T}_h}, \quad (2.2b)$$

$$\langle \hat{\mathbf{q}}_h \cdot \mathbf{n}, \mu \rangle_{\partial\mathcal{T}_h \setminus \partial\Omega} = 0, \quad (2.2c)$$

$$\hat{u}_h = 0 \quad \text{on } \partial\Omega. \quad (2.2d)$$

with the normal component of the numerical trace defined

$$\hat{\mathbf{q}}_h \cdot \mathbf{n} = \mathbf{q}_h \cdot \mathbf{n} + \tau(u_h - \hat{u}_h) \quad \text{on } \partial\mathcal{T}_h, \quad (2.2e)$$

where τ is a local stabilization parameter.

The Eqs. (2.2a) and (2.2b) can be derived by combining the numerical trace \hat{u}_h and the numerical flux $\hat{\mathbf{q}}_h$ from the result of integration by parts in the Eqs. (2.1a) and (2.1b), respectively. Numerical trace \hat{u}_h and numerical flux $\hat{\mathbf{q}}_h$ are approximations of u and \mathbf{q} on the element interface $\partial\mathcal{T}_h$, respectively. Also, numerical trace is used to solve the global problem.

For the convenience of analysis, we define the weighted L^2 -norm with a local stabilization parameter τ as following:

$$\|w\|_{\tau, \partial\mathcal{T}_h}^2 := \langle \tau w, w \rangle_{\partial\mathcal{T}_h}.$$

We will use the HDG projection Π_h to show the well-posedness and error estimation. For this, we will introduce the definition of Π_h and the preliminary result of [10] without the proof. The projection Π_h into $\mathbf{V}_h \times W_h$ is defined as follows.

Given $(u, \mathbf{q}) \in H^1(\mathcal{T}_h) \times \mathbf{H}_{div}(\mathcal{T}_h)$, the function $\Pi_h(u, \mathbf{q}) = (\Pi_W u, \Pi_V \mathbf{q})$ on an arbitrary simplex $T \in \mathcal{T}_h$ is the element of $W_h \times \mathbf{V}_h$ which solves

$$(\Pi_W u, w)_T = (u, w)_T, \quad \forall w \in \mathcal{P}_{k-1}(T), \quad (2.3a)$$

$$(\Pi_V \mathbf{q}, \mathbf{v})_T = (\mathbf{q}, \mathbf{v})_T, \quad \forall \mathbf{v} \in \mathcal{P}_{k-1}(T) \quad (2.3b)$$

$$\langle \Pi_V \mathbf{q} \cdot \mathbf{n} + \tau \Pi_W u, \mu \rangle_F = \langle \mathbf{q} \cdot \mathbf{n} + \tau u, \mu \rangle_F, \quad \forall \mu \in \mathcal{P}_k(F), \quad (2.3c)$$

for all interfaces F of the simplex T . Also, P_M denotes the L^2 -orthogonal projection onto M_h . From the last Eq. (2.3c) and the definition of the projection operator P_M , we immediately have

$$P_M(\mathbf{q} \cdot \mathbf{n}) + \tau P_M u = \Pi_V \mathbf{q} \cdot \mathbf{n} + \tau \Pi_W u, \quad \text{for all } F \in \partial\mathcal{T}_h. \quad (2.4)$$

The following result show that the HDG projections are well defined and give approximation properties. In our analysis, we will derive the error estimations for $\|u - u_h\|_\Omega$ and $\|\mathbf{q} - \mathbf{q}_h\|_\Omega$ by combining the result of Lemma 2.1 with results derived from error Eqs. (3.6).

Lemma 2.1. *If the local spaces are given by $(W(T), \mathbf{V}(T)) = (\mathcal{P}_k(T), \mathcal{P}_k(T))$ for $k \geq 0$ and τ is nonnegative, then the system (2.3) is uniquely solvable for $\Pi_W u$ and $\Pi_V \mathbf{q}$. Furthermore, there is a constant C independent from the choice of $T \in \mathcal{T}_h$ and τ such that for all $(u, \mathbf{q}) \in H^s(T) \times \mathbf{H}^s(T)$ and $1 \leq s \leq k + 1$,*

$$\begin{aligned} \|\mathbf{q} - \Pi_V \mathbf{q}\|_T &\leq Ch^s (\|\mathbf{q}\|_{s,T} + \tau \|u\|_{s,T}), \\ \|u - \Pi_W u\|_T &\leq Ch^s (\|u\|_{s,T} + \tau^{-1} \|\nabla \cdot \mathbf{q}\|_{s-1,T}). \end{aligned}$$

3. ANALYSIS

3.1. Assumptions. We will provide necessary assumptions to analyze the HDG method for elliptic problems with nonlinear coefficients. First of all, we give the assumption related to nonlinear coefficients $\alpha(u)$. Recall that $\alpha(u) = \kappa(u)^{-1}$.

Assumption 3.1. *α is chosen in such a way that there exist positive constants α_1, α_2 and L such that for all $u, v \in \mathbb{R}$, we have the following inequalities:*

- i) $0 < \alpha_1 \leq \alpha(u) \leq \alpha_2 < \infty$,
- ii) $|\alpha(u) - \alpha(v)| \leq L \|u - v\|_\Omega$.

To ensure the existence and uniqueness of the solution and to derive the error estimations in the HDG method, we consider the following dual problem for any given $\Psi \in L^2(\Omega)$:

$$\boldsymbol{\theta} + \nabla \phi = 0, \quad \text{in } \Omega \tag{3.1a}$$

$$\nabla \cdot \boldsymbol{\theta} = \Psi, \quad \text{in } \Omega \tag{3.1b}$$

$$\phi = 0, \quad \text{on } \partial\Omega. \tag{3.1c}$$

Assumption 3.2. *We assume that the dual problem (3.1) admits the H^2 -regularity*

$$\|\phi\|_{2,\Omega} + \|\boldsymbol{\theta}\|_{1,\Omega} \leq C_{reg} \|\Psi\|_\Omega, \tag{3.2}$$

for all $\Psi \in L^2(\Omega)$.

We consider the dual problem (3.1) without the nonlinear coefficient α . Notice that the Assumption 3.2 may not be valid if the coefficient α is not smooth or with high contrast. If Ω is convex polygon, the above assumption holds in [28].

In the HDG method, the well-posedness and accuracy involves a local stabilization parameter τ in [2, 10]. The existence and uniqueness of the numerical solution for linear elliptic problems should be guaranteed by the following assumption. This assumption shall be satisfied with our problems.

Assumption 3.3. *We assume that a local stabilization parameter τ is a positive constant.*

3.2. Well-posedness. Now we prove the well-posedness of the nonlinear problem (2.2) based on the Banach fixed-point theorem in [29]. We define an operator $\mathcal{O} : W_h \rightarrow W_h$ mapping η_h to u_h . Here, u_h is the first component of the approximation $(u_h, \mathbf{q}_h, \hat{u}_h) \in W_h \times \mathbf{V}_h \times M_h$ satisfying the following systems:

$$(\alpha(\eta_h)\mathbf{q}_h, \mathbf{v})_{\mathcal{T}_h} - (u_h, \nabla \cdot \mathbf{v})_{\mathcal{T}_h} + \langle \hat{u}_h, \mathbf{v} \cdot \mathbf{n} \rangle_{\partial\mathcal{T}_h} = 0, \quad (3.3a)$$

$$-(\mathbf{q}_h, \nabla w)_{\mathcal{T}_h} + \langle \hat{\mathbf{q}}_h \cdot \mathbf{n}, w \rangle_{\partial\mathcal{T}_h} = (f, w)_{\mathcal{T}_h}, \quad (3.3b)$$

$$\langle \hat{\mathbf{q}}_h \cdot \mathbf{n}, \mu \rangle_{\partial\mathcal{T}_h} = 0 \quad (3.3c)$$

$$\hat{u}_h = 0 \quad \text{on} \quad \partial\Omega, \quad (3.3d)$$

for all $(w, \mathbf{v}, \mu) \in W_h \times \mathbf{V}_h \times M_h$, with a numerical trace for the flux defined

$$\hat{\mathbf{q}}_h \cdot \mathbf{n} = \mathbf{q}_h \cdot \mathbf{n} + \tau(u_h - \hat{u}_h) \quad \text{on} \quad \partial\mathcal{T}_h. \quad (3.3e)$$

For a given $\eta_h \in V_h$, the system (3.3) is the HDG formulation for the linear elliptic equation. The existence and uniqueness of the numerical solution for the linear systems have been shown in [2]. So, the mapping \mathcal{O} is well-defined. If the mapping \mathcal{O} is a contraction, that is, $\|\mathcal{O}(\eta_h^1) - \mathcal{O}(\eta_h^2)\|_{\Omega} < \|\eta_h^1 - \eta_h^2\|_{\Omega}$, the discrete problem (2.2) is well-posed. We will use the properties listed in the lemmas below to show that this inequality holds.

Lemma 3.4. *We assume that the Assumptions 3.1 and 3.3 are satisfied. Then we have*

$$\|\alpha(\eta_h)^{\frac{1}{2}}\mathbf{q}_h\|_{\Omega}^2 + \|u_h - \hat{u}_h\|_{\tau, \partial\mathcal{T}_h}^2 \leq \|f\|_{\Omega}\|u_h\|_{\Omega}. \quad (3.4)$$

Proof. Take $(w, \mathbf{v}, \mu) = (u_h, \mathbf{q}_h, \hat{u}_h)$ in Eqs. (3.3a)-(3.3c). Adding, we get, after some algebraic manipulation,

$$(\alpha(\eta_h)\mathbf{q}_h, \mathbf{q}_h)_{\mathcal{T}_h} + \Gamma_h = (f, u_h)_{\mathcal{T}_h}$$

where

$$\Gamma_h = -(u_h, \nabla \cdot \mathbf{q}_h)_{\mathcal{T}_h} + \langle \hat{u}_h, \mathbf{q}_h \cdot \mathbf{n} \rangle_{\partial\mathcal{T}_h} - (\mathbf{q}_h, \nabla u_h)_{\mathcal{T}_h} + \langle \hat{\mathbf{q}}_h \cdot \mathbf{n}, u_h - \hat{u}_h \rangle_{\partial\mathcal{T}_h}.$$

By integrating by parts and using the Eq. (3.3e), we get

$$\begin{aligned} \Gamma_h &= \langle \hat{u}_h, \mathbf{q}_h \cdot \mathbf{n} \rangle_{\partial\mathcal{T}_h} - \langle u_h, \mathbf{q}_h \cdot \mathbf{n} \rangle_{\partial\mathcal{T}_h} + \langle \hat{\mathbf{q}}_h \cdot \mathbf{n}, u_h - \hat{u}_h \rangle_{\partial\mathcal{T}_h} \\ &= \langle \hat{\mathbf{q}}_h \cdot \mathbf{n}, u_h - \hat{u}_h \rangle_{\partial\mathcal{T}_h} - \langle \mathbf{q}_h \cdot \mathbf{n}, u_h - \hat{u}_h \rangle_{\partial\mathcal{T}_h} \\ &= \langle \hat{\mathbf{q}}_h \cdot \mathbf{n} - \mathbf{q}_h \cdot \mathbf{n}, u_h - \hat{u}_h \rangle_{\partial\mathcal{T}_h} \\ &= \langle \tau(u_h - \hat{u}_h), u_h - \hat{u}_h \rangle_{\partial\mathcal{T}_h} = \|u_h - \hat{u}_h\|_{\tau, \partial\mathcal{T}_h}^2. \end{aligned}$$

Applying Cauchy-Schwarz inequality, we have

$$\|\alpha(\eta_h)^{\frac{1}{2}}\mathbf{q}_h\|_{\Omega}^2 + \|u_h - \hat{u}_h\|_{\tau, \partial\mathcal{T}_h}^2 \leq \|f\|_{\Omega}\|u_h\|_{\Omega}.$$

This completes the proof. \square

We will use of the properties of the HDG projections (Π_W, Π_V) and the dual system (3.1) with $\Psi = u_h \in L^2(\Omega)$.

Lemma 3.5. *We assume that Assumptions 3.1, 3.2, and 3.3 are satisfied. If $\tau = \alpha_1^{-1}$, then we have*

$$\|u_h\|_{\Omega} \leq C_{reg}^2 \|\alpha(\eta_h)\| \|f\|_{\Omega}. \quad (3.5)$$

Proof. Let ϕ and $\boldsymbol{\theta}$ be the solutions to the dual problem (3.1) with $\Psi = u_h \in L^2(\Omega)$. We begin by using the Eq. (3.1b) to write that

$$\begin{aligned} \|u_h\|_{\Omega}^2 &= (u_h, u_h)_{\mathcal{T}_h} = (u_h, \nabla \cdot \boldsymbol{\theta})_{\mathcal{T}_h} \\ &= (u_h, \nabla \cdot \mathbf{\Pi}_V \boldsymbol{\theta})_{\mathcal{T}_h} + (u_h, \nabla \cdot (\boldsymbol{\theta} - \mathbf{\Pi}_V \boldsymbol{\theta}))_{\mathcal{T}_h}. \end{aligned}$$

By integrating by parts for the second term of the above equation and using the property of $\mathbf{\Pi}_V$, we get

$$\begin{aligned} \|u_h\|_{\Omega}^2 &= (u_h, \nabla \cdot \mathbf{\Pi}_V \boldsymbol{\theta})_{\mathcal{T}_h} - (\nabla u_h, \boldsymbol{\theta} - \mathbf{\Pi}_V \boldsymbol{\theta})_{\mathcal{T}_h} + \langle u_h, (\boldsymbol{\theta} - \mathbf{\Pi}_V \boldsymbol{\theta}) \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} \\ &= (u_h, \nabla \cdot \mathbf{\Pi}_V \boldsymbol{\theta})_{\mathcal{T}_h} + \langle u_h, (\boldsymbol{\theta} - \mathbf{\Pi}_V \boldsymbol{\theta}) \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} \end{aligned}$$

Taking $\mathbf{v} = \mathbf{\Pi}_V \boldsymbol{\theta}$ in the Eq. (3.3a), we observe that

$$(u_h, \nabla \cdot \mathbf{\Pi}_V \boldsymbol{\theta})_{\mathcal{T}_h} = (\alpha(\eta_h) \mathbf{q}_h, \mathbf{\Pi}_V \boldsymbol{\theta})_{\mathcal{T}_h} + \langle \widehat{u}_h, \mathbf{\Pi}_V \boldsymbol{\theta} \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h}.$$

Now using the fact that \widehat{u}_h are single-valued on $\partial \mathcal{T}_h$, so $\langle \widehat{u}_h, \boldsymbol{\theta} \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} = 0$ after some algebraic manipulation, we obtain

$$\|u_h\|_{\Omega}^2 = (\alpha(\eta_h) \mathbf{q}_h, \mathbf{\Pi}_V \boldsymbol{\theta})_{\mathcal{T}_h} + \langle u_h - \widehat{u}_h, (\boldsymbol{\theta} - \mathbf{\Pi}_V \boldsymbol{\theta}) \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h}.$$

We now estimate the above two terms separately. By Assumptions 3.1 and 3.2, we get

$$\begin{aligned} |(\alpha(\eta_h) \mathbf{q}_h, \mathbf{\Pi}_V \boldsymbol{\theta})_{\mathcal{T}_h}| &\leq \alpha(\eta_h)^{\frac{1}{2}} \left| (\alpha(\eta_h)^{\frac{1}{2}} \mathbf{q}_h, \mathbf{\Pi}_V \boldsymbol{\theta})_{\mathcal{T}_h} \right| \\ &\leq \alpha(\eta_h)^{\frac{1}{2}} \left(\left| (\alpha(\eta_h)^{\frac{1}{2}} \mathbf{q}_h, \boldsymbol{\theta})_{\mathcal{T}_h} \right| + \left| (\alpha(\eta_h)^{\frac{1}{2}} \mathbf{q}_h, \mathbf{\Pi}_V \boldsymbol{\theta} - \boldsymbol{\theta})_{\mathcal{T}_h} \right| \right) \\ &\leq \alpha(\eta_h)^{\frac{1}{2}} \|\alpha(\eta_h)^{\frac{1}{2}} \mathbf{q}_h\|_{\Omega} (\|\boldsymbol{\theta}\|_{\Omega} + \|\mathbf{\Pi}_V \boldsymbol{\theta} - \boldsymbol{\theta}\|_{\Omega}) \\ &\leq \alpha(\eta_h)^{\frac{1}{2}} \|\alpha(\eta_h)^{\frac{1}{2}} \mathbf{q}_h\|_{\Omega} (\|\boldsymbol{\theta}\|_{\Omega} + h\|\phi\|_{1,\Omega}) \quad \text{by } h \ll 1 \\ &\leq C_{reg} \alpha(\eta_h)^{\frac{1}{2}} \|\alpha(\eta_h)^{\frac{1}{2}} \mathbf{q}_h\|_{\Omega} \|u_h\|_{\Omega}. \end{aligned}$$

Similarly,

$$\begin{aligned} |\langle u_h - \widehat{u}_h, (\boldsymbol{\theta} - \mathbf{\Pi}_V \boldsymbol{\theta}) \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h}| &\leq \tau^{-\frac{1}{2}} h^{-\frac{1}{2}} \|\| u_h - \widehat{u}_h \|\|_{\tau, \partial \mathcal{T}_h} \|\boldsymbol{\theta} - \mathbf{\Pi}_V \boldsymbol{\theta}\|_{\Omega} \\ &\leq C_{reg} \tau^{-\frac{1}{2}} h^{\frac{1}{2}} \|\| u_h - \widehat{u}_h \|\|_{\tau, \partial \mathcal{T}_h} \|u_h\|_{\Omega} \\ &\leq C_{reg} \tau^{-\frac{1}{2}} \|\| u_h - \widehat{u}_h \|\|_{\tau, \partial \mathcal{T}_h} \|u_h\|_{\Omega} \quad \text{by } h \ll 1. \end{aligned}$$

Then, applying Lemma 3.4, we get

$$\begin{aligned} \|u_h\|_{\Omega} &\leq C_{reg} \max\{\alpha(\eta_h)^{\frac{1}{2}}, \tau^{-\frac{1}{2}}\} \left(\|\alpha(\eta_h)^{\frac{1}{2}} \mathbf{q}_h\|_{\Omega} + \|\| u_h - \widehat{u}_h \|\|_{\tau, \partial \mathcal{T}_h} \right) \\ &\leq C_{reg} \max\{\alpha(\eta_h)^{\frac{1}{2}}, \tau^{-\frac{1}{2}}\} \|f\|_{\Omega}^{\frac{1}{2}} \|u_h\|_{\Omega}^{\frac{1}{2}}. \end{aligned}$$

Since $\tau^{-1} \leq \alpha(\eta_h)$, we have

$$\|u_h\|_{\Omega} \leq C_{reg}^2 |\alpha(\eta_h)| \|f\|_{\Omega}.$$

□

Now we will show that the operator \mathcal{O} is a contraction mapping. By combining with the Banach fixed point theorem, the following theorem guarantees that the discrete system has a unique solution.

Theorem 3.6. *Suppose that Assumptions 3.1, 3.2, and 3.3 are satisfied. If $\tau = a_1^{-1}$ and $C_{reg}^2 L \|f\|_\Omega < 1$, then \mathcal{O} is a contraction operator.*

Proof. Let $\eta_h^1, \eta_h^2 \in W_h$ and set $u_h^1 := \mathcal{O}(\eta_h^1)$, $u_h^2 := \mathcal{O}(\eta_h^2)$. Then u_h^1 and u_h^2 are solutions of the system (3.3). By Lemma 3.5, we have

$$\begin{aligned} \|\mathcal{O}(\eta_h^1) - \mathcal{O}(\eta_h^2)\|_\Omega &= \|u_h^1 - u_h^2\|_\Omega \\ &\leq C_{reg}^2 |\alpha(\eta_h^1) - \alpha(\eta_h^2)| \|f\|_\Omega \leq C_{reg}^2 L \|\eta_h^1 - \eta_h^2\|_\Omega \|f\|_\Omega \\ &< \|\eta_h^1 - \eta_h^2\|_\Omega. \end{aligned}$$

Therefore, \mathcal{O} is a contraction operator. This completes the proof. \square

3.3. Error estimations. To derive the error estimations for the system (4.2), we define the projection errors as follows:

$$\begin{aligned} e_u &:= \Pi_W u - u_h, & e_q &:= \Pi_V \mathbf{q} - \mathbf{q}_h, \\ e_{\hat{u}} &:= P_M u - \hat{u}_h, & e_{\hat{q}} \cdot \mathbf{n} &:= P_M(\mathbf{q} \cdot \mathbf{n}) - \hat{\mathbf{q}}_h \cdot \mathbf{n}, \\ \delta_u &:= u - \Pi_W u, & \delta_q &:= \mathbf{q} - \Pi_V \mathbf{q}, \end{aligned}$$

where (u, \mathbf{q}) and $(u_h, \mathbf{q}_h, \hat{u}_h)$ were solutions of the systems (2) and (5), respectively.

We note that by triangle inequality,

$$\|u_h - u\|_\Omega \leq \|u - \Pi_W u\|_\Omega + \|\Pi_W u - u_h\|_\Omega.$$

The first term (i.e., $\|u - \Pi_W u\|_\Omega$) that appears in RHS is bounded by Lemma 2.1. Hence, we only need to find an upper bound for $\Pi_W u - u_h =: e_u$. Similarly, we will find an upper bound for e_q to get the error estimation for $\|\mathbf{q}_h - \mathbf{q}\|_\Omega$.

We begin by obtaining the error equations to find an upper bound for e_u and e_q . The proofs follow the technique developed in [1, 10]

Lemma 3.7. *We have*

$$\begin{aligned} (\alpha(u_h) e_q, \mathbf{v})_{\mathcal{T}_h} - (e_u, \nabla \cdot \mathbf{v})_{\mathcal{T}_h} + \langle e_{\hat{u}}, \mathbf{v} \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} &= -(\alpha(u_h) \delta_q, \mathbf{v})_{\mathcal{T}_h} \\ &\quad + ((\alpha(u_h) - \alpha(u)) \mathbf{q}, \mathbf{v})_{\mathcal{T}_h}, \end{aligned} \quad (3.6a)$$

$$-(e_q, \nabla w)_{\mathcal{T}_h} + \langle e_{\hat{q}} \cdot \mathbf{n}, w \rangle_{\partial \mathcal{T}_h} = 0, \quad (3.6b)$$

$$\langle e_{\hat{q}} \cdot \mathbf{n}, \mu \rangle_{\partial \mathcal{T}_h \setminus \partial \Omega} = 0, \quad (3.6c)$$

for all $(w, \mathbf{v}, \mu) \in W_h \times \mathbf{V}_h \times M_h$, where

$$e_{\hat{q}} \cdot \mathbf{n} = e_q \cdot \mathbf{n} + \tau(e_u - e_{\hat{u}}) \quad \text{on} \quad \partial \mathcal{T}_h. \quad (3.6d)$$

Proof. The exact solution (u, \mathbf{q}) obviously satisfies

$$\begin{aligned} (\alpha(u) \mathbf{q}, \mathbf{v})_{\mathcal{T}_h} - (u, \nabla \cdot \mathbf{v})_{\mathcal{T}_h} + \langle u, \mathbf{v} \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} &= 0, \\ -(\mathbf{q}, \nabla w)_{\mathcal{T}_h} + \langle \mathbf{q} \cdot \mathbf{n}, w \rangle_{\partial \mathcal{T}_h} &= (f, w), \\ \langle \mathbf{q} \cdot \mathbf{n}, \mu \rangle_{\partial \mathcal{T}_h \setminus \partial \Omega} &= 0, \end{aligned}$$

for all $(w, \mathbf{v}, \mu) \in W_h \times \mathbf{V}_h \times M_h$.

By the definition of the projections (Π_W, Π_V, P_M) , we can get

$$\begin{aligned} (\alpha(u_h))\Pi_V \mathbf{q}, \mathbf{v})_{\mathcal{T}_h} - (\Pi_W u, \nabla \cdot \mathbf{v})_{\mathcal{T}_h} + \langle P_M u, \mathbf{v} \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h} &= (\alpha(u_h)\Pi_V \mathbf{q}, \mathbf{v})_{\mathcal{T}_h} \\ &\quad - (\alpha(u)\mathbf{q}, \mathbf{v})_{\mathcal{T}_h}, \\ -(\Pi_V \mathbf{q}, \nabla w)_{\mathcal{T}_h} + \langle P_M(\mathbf{q} \cdot \mathbf{n}), w \rangle_{\partial \mathcal{T}_h} &= (f, w), \\ \langle P_M(\mathbf{q} \cdot \mathbf{n}), \mu \rangle_{\partial \mathcal{T}_h \setminus \partial \Omega} &= 0, \end{aligned}$$

for all $(w, \mathbf{v}, \mu) \in W_h \times \mathbf{V}_h \times M_h$. Subtracting the first three equations defining the HDG method (i.e., (2.2a)-(2.2c)), from the above equations in order, we obtain (3.6a)-(3.6c).

It remains to prove the identity (3.6d) for $e_{\hat{q}} \cdot \mathbf{n}$. On each interface $F \in \partial T$, by using the definition of numerical trace (2.2e), we have

$$\begin{aligned} e_{\hat{q}} \cdot \mathbf{n} - e_q \cdot \mathbf{n} &= P_M(\mathbf{q} \cdot \mathbf{n}) - \hat{\mathbf{q}}_h \cdot \mathbf{n} - (\Pi_V \mathbf{q} \cdot \mathbf{n} - \mathbf{q}_h \cdot \mathbf{n}) \\ &= P_M(\mathbf{q} \cdot \mathbf{n}) - \Pi_V \mathbf{q} \cdot \mathbf{n} - (\hat{\mathbf{q}}_h \cdot \mathbf{n} - \mathbf{q}_h \cdot \mathbf{n}) \\ &= P_M(\mathbf{q} \cdot \mathbf{n}) - \Pi_V \mathbf{q} \cdot \mathbf{n} - \tau(u_h - \hat{u}_h). \end{aligned}$$

Then using the property (2.4) of the projections Π_V and Π_W , the equality reduces to

$$\begin{aligned} e_{\hat{q}} \cdot \mathbf{n} - e_q \cdot \mathbf{n} &= \tau(-P_M u + \Pi_W u) - \tau(u_h - \hat{u}_h) \\ &= \tau(e_u - e_{\hat{u}}). \end{aligned}$$

This completes the proof. \square

Lemma 3.8 shows the error estimation for e_q .

Lemma 3.8. *We assume that Assumption 3.1 and 3.3 are satisfied, then we have*

$$\|e_q\|_{\Omega}^2 + \|e_u - e_{\hat{u}}\|_{\tau, \partial \mathcal{T}_h}^2 \leq C (\|\delta_q\|_{\Omega}^2 + \|\delta_u\|_{\Omega}^2) + C_L \|e_u\|_{\Omega}^2 \quad (3.7)$$

where C_L is dependent on the Lipschitz constant L .

Proof. Take $(w, \mathbf{v}, \mu) = (e_u, e_q, \mathbf{v})$ in the error equations (3.6a)-(3.6c). Similar to the proof of Lemma 3.4, we have the following identity

$$(\alpha(u_h)e_q, e_q)_{\mathcal{T}_h} + \|e_u - e_{\hat{u}}\|_{\tau, \partial \mathcal{T}_h}^2 = \mathbb{T}_1 + \mathbb{T}_2,$$

where

$$\mathbb{T}_1 = -(\alpha(u_h)\delta_q, e_q)_{\mathcal{T}_h} \quad \text{and} \quad \mathbb{T}_2 = ((\alpha(u_h) - \alpha(u))\mathbf{q}, e_q)_{\mathcal{T}_h}.$$

We estimate the above terms separately. Applying Cauchy-Schwarz and Young's inequalities, we get

$$|\mathbb{T}_1| = |(\alpha(u_h)\delta_q, e_q)_{\mathcal{T}_h}| \leq \alpha_2 |(\delta_q, e_q)_{\mathcal{T}_h}| \leq \frac{\alpha_2}{\alpha_1} \|\delta_q\|_{\Omega}^2 + \frac{\alpha_1}{4} \|e_q\|_{\Omega}^2.$$

Since Ω is a bounded domain and the weak solution of the elliptic equation lies in various higher Sobolev spaces [30], we set $\|\mathbf{q}\|_{\infty} = C$. Then, we have

$$\begin{aligned} |\mathbb{T}_2| &= |((\alpha(u_h) - \alpha(u))\mathbf{q}, e_q)_{\mathcal{T}_h}| \leq CL \|u_h - u\|_{\Omega} \|e_q\|_{\Omega} \\ &\leq CL (\|e_u\|_{\Omega} + \|\delta_u\|_{\Omega}) \|e_q\|_{\Omega} \\ &\leq \frac{CL}{\alpha_1} (\|e_u\|_{\Omega}^2 + \|\delta_u\|_{\Omega}^2) + \frac{\alpha_1}{4} \|e_q\|_{\Omega}^2. \end{aligned}$$

Since $\alpha_1 \|e_q\|_\Omega^2 \leq (\alpha(u_h)e_q, e_q)_{\mathcal{T}_h}$, we get

$$\frac{\alpha_1}{2} \|e_q\|_\Omega^2 + \|e_u - e_{\hat{u}}\|_{\tau, \partial\mathcal{T}_h}^2 \leq \frac{\alpha_2}{\alpha_1} \|\delta_q\|_\Omega^2 + \frac{CL}{\alpha_1} \|\delta_u\|_\Omega^2 + \frac{CL}{\alpha_1} \|e_u\|_\Omega^2.$$

Therefore, we have

$$\|e_q\|_\Omega^2 + \|e_u - e_{\hat{u}}\|_{\tau, \partial\mathcal{T}_h}^2 \leq C_1 (\|\delta_q\|_\Omega^2 + \|\delta_u\|_\Omega^2) + C_L \|e_u\|_\Omega^2,$$

where

$$C_1 = \max \left\{ \frac{\alpha_2}{\alpha_1 \min\{\frac{\alpha_1}{2}, 1\}}, C_L \right\} \quad \text{with} \quad C_L = \frac{CL}{\alpha_1 \min\{\frac{\alpha_1}{2}, 1\}}.$$

This completes the proof. \square

Next, we derive the error estimation for e_u .

Lemma 3.9. *We assume that Assumptions 3.1, 3.2, and 3.3 are satisfied, and $L < 1/C_{reg}\|\mathbf{q}\|_\infty$ holds with the Lipschitz constant L , then we have*

$$\|e_u\|_\Omega \leq C (\|\delta_u\|_\Omega + \|\delta_q\|_\Omega) + C_L^* (\|e_q\|_\Omega + \|e_u - e_{\hat{u}}\|_{\tau, \partial\mathcal{T}_h}), \quad (3.8)$$

where C_L^* is dependent on the Lipschitz constant L .

Proof. Let ϕ and θ be the solutions to the dual problem (3.1) with $\Psi = e_u \in L^2(\Omega)$. Similar to the proof of Lemma 3.5, we have the following identity

$$\|e_u\|_\Omega^2 = \mathbb{S}_1 + \mathbb{S}_2 + \mathbb{S}_3,$$

where

$$\begin{aligned} \mathbb{S}_1 &= ((\alpha(u) - \alpha(u_h))\mathbf{q}, \mathbf{\Pi}_V\theta)_{\mathcal{T}_h}, & \mathbb{S}_2 &= (\alpha(u_h)(\mathbf{q} - \mathbf{q}_h), \mathbf{\Pi}_V\theta)_{\mathcal{T}_h}, \\ \mathbb{S}_3 &= \langle e_u - e_{\hat{u}}, (\theta - \mathbf{\Pi}_V\theta) \cdot \mathbf{n} \rangle_{\partial\mathcal{T}_h}. \end{aligned}$$

We estimate the above three terms separately. By Assumptions 3.1 and 3.2, and $\|\mathbf{q}\|_\infty = C$, we have

$$\begin{aligned} |\mathbb{S}_1| &\leq |((\alpha(u) - \alpha(u_h))\mathbf{q}, \theta)_{\mathcal{T}_h}| + |((\alpha(u) - \alpha(u_h))\mathbf{q}, \mathbf{\Pi}_V\theta - \theta)_{\mathcal{T}_h}| \\ &\leq CL\|u - u_h\|_\Omega (\|\theta\|_\Omega + \|\mathbf{\Pi}_V\theta - \theta\|_\Omega) \\ &\leq CL(\|e_u\|_\Omega + \|\delta_u\|_\Omega) (\|\theta\|_\Omega + \|\mathbf{\Pi}_V\theta - \theta\|_\Omega) \\ &\leq CL(\|e_u\|_\Omega + \|\delta_u\|_\Omega) (\|\theta\|_\Omega + h\|\phi\|_{1,\Omega}) \quad \text{by } h \ll 1 \\ &\leq C_{reg}CL(\|e_u\|_\Omega + \|\delta_u\|_\Omega) \|e_u\|_\Omega. \end{aligned}$$

Similarly,

$$\begin{aligned} |\mathbb{S}_2| &\leq |(\alpha(u_h)(\mathbf{q} - \mathbf{q}_h), \theta)_{\mathcal{T}_h}| + |(\alpha(u_h)(\mathbf{q} - \mathbf{q}_h), \mathbf{\Pi}_V\theta - \theta)_{\mathcal{T}_h}| \\ &\leq \alpha_2 \|\mathbf{q} - \mathbf{q}_h\|_\Omega (\|\theta\|_\Omega + \|\mathbf{\Pi}_V\theta - \theta\|_\Omega) \\ &\leq \alpha_2 (\|e_q\|_\Omega + \|\delta_q\|_\Omega) (\|\theta\|_\Omega + \|\mathbf{\Pi}_V\theta - \theta\|_\Omega) \\ &\leq C_{reg}\alpha_2 (\|e_q\|_\Omega + \|\delta_q\|_\Omega) \|e_u\|_\Omega. \end{aligned}$$

We apply Cauchy-Schwartz inequality to obtain

$$\begin{aligned}
|\mathbb{S}_3| &= |\langle e_u - e_{\hat{u}}, (\boldsymbol{\theta} - \mathbf{\Pi}_V \boldsymbol{\theta}) \cdot \mathbf{n} \rangle_{\partial\mathcal{T}_h}| \\
&\leq \tau^{-\frac{1}{2}} \| \| e_u - e_{\hat{u}} \| \|_{\tau, \partial\mathcal{T}_h} \| (\boldsymbol{\theta} - \mathbf{\Pi}_V \boldsymbol{\theta}) \cdot \mathbf{n} \|_{\partial\mathcal{T}_h} \\
&\leq C_{reg} \tau^{-\frac{1}{2}} h^{\frac{1}{2}} \| \| e_u - e_{\hat{u}} \| \|_{\tau, \partial\mathcal{T}_h} \| e_u \|_{\Omega} \\
&\leq C_{reg} \tau^{-\frac{1}{2}} \| \| e_u - e_{\hat{u}} \| \|_{\tau, \partial\mathcal{T}_h} \| e_u \|_{\Omega} \quad \text{by } h \ll 1.
\end{aligned}$$

Since $C_{reg}CL < 1$, we get

$$\begin{aligned}
(1 - C_{reg}CL) \| e_u \| &\leq C_{reg} \max\{CL, \alpha_2\} (\| \delta_u \|_{\Omega} + \| \boldsymbol{\delta}_q \|_{\Omega}) \\
&\quad + C_{reg} \max\{\alpha_2, \tau^{-\frac{1}{2}}\} (\| \mathbf{e}_q \|_{\Omega} + \| \| e_u - e_{\hat{u}} \| \|_{\tau, \partial\mathcal{T}_h}).
\end{aligned}$$

Therefore, we have

$$\| e_u \|_{\Omega} \leq C_2 (\| \delta_u \|_{\Omega} + \| \boldsymbol{\delta}_q \|_{\Omega}) + C_L^* (\| \mathbf{e}_q \|_{\Omega} + \| \| e_u - e_{\hat{u}} \| \|_{\tau, \partial\mathcal{T}_h}),$$

where

$$C_2 = \frac{C_{reg} \max\{CL, \alpha_2\}}{1 - C_{reg}CL} \quad \text{and} \quad C_L^* = \frac{C_{reg} \max\{\alpha_2, \tau^{-\frac{1}{2}}\}}{1 - C_{reg}CL}.$$

This completes the proof. \square

Lemma 3.8 and 3.9 show the upper bound of the projection errors for \mathbf{e}_q and e_u , respectively. However, the upper bound of the error \mathbf{e}_q contains e_u , and the upper bound of the error e_u contains \mathbf{e}_q . To overcome this, we need to consider an additional condition. The following theorem shows that under the additional condition, each projection error has an optimal order of convergence.

Theorem 3.10. *Let the condition of Lemma 3.9 be satisfied. Let L , C_L , and C_L^* be the Lipschitz constant, the constant of Lemma 3.8, and the constant of Lemma 3.9, respectively. If the condition $C_L C_L^* < 1$, then for all $0 \leq s \leq k + 1$*

$$\| \mathbf{e}_q \|_{\Omega} + \| \| e_u - e_{\hat{u}} \| \|_{\tau, \partial\mathcal{T}_h} \leq Ch^s (\| \mathbf{q} \|_s + \| u \|_s) \quad (3.9)$$

and

$$\| e_u \|_{\Omega} \leq Ch^s (\| \mathbf{q} \|_s + \| u \|_s), \quad (3.10)$$

where C depends on the stabilization parameter τ .

Proof. It follows from Lemma 3.8 and 3.9, that

$$\begin{aligned}
\| \mathbf{e}_q \|_{\Omega} + \| \| e_u - e_{\hat{u}} \| \|_{\tau, \partial\mathcal{T}_h} &\leq C (\| \boldsymbol{\delta}_q \|_{\Omega} + \| \delta_u \|_{\Omega}) + L^* \| e_u \|_{\Omega} \\
&\leq C (\| \boldsymbol{\delta}_q \|_{\Omega} + \| \delta_u \|_{\Omega}) + C_L C_L^* (\| \mathbf{e}_q \|_{\Omega} + \| \| e_u - e_{\hat{u}} \| \|_{\tau, \partial\mathcal{T}_h}).
\end{aligned}$$

By $C_L C_L^* < 1$ and Lemma 2.1, we get the estimate (3.9). Similarly, we get the estimate (3.10). This completes the proof. \square

Remark 3.11. *If the Lipschitz constant L is small enough, the additional conditions of Lemma 3.8 and Lemma 3.9 are satisfied. Since the value of the Lipschitz constant depends on the nonlinear coefficient α (or κ), if the nonlinear term has good properties, we can obtain the HDG approximation with the same convergence ratio as in linear problems.*

Remark 3.12. *The constant C in Theorem 3.10 depends on the value of the local stabilization parameter τ according to the estimates of Lemma 2.1 and 3.9. To emphasize the optimal convergence order, we express the constant C depending on τ . However, if $\tau = 1$, or if τ is a fixed constant, then both approximation errors converge with the optimal order of $k + 1$, when the functions u and \mathbf{q} are smooth enough.*

4. NUMERICAL RESULTS

4.1. Matrix formulation for iterative techniques. For implementation, we need to use an iterative method to find the solution of nonlinear problems with linear structure. So, we consider the linear problems (3.3). We can consider a variety of iterative techniques. The basic concept is to update the solution at the current iteration step using the previous solution. Also, we can find a final solution using a suitable tolerance between the current and previous solutions.

We introduce the iterative solver. Let k be the iteration count and $\alpha(\eta_h) = \alpha(u_h^{k-1})$. Then we can get the solution $(u_h^k, \mathbf{q}_h^k, \widehat{u}_h^k)$ of the system (3.3). Also, for each element T , we can find the pair $(u_h^k(\widehat{u}_h^k, f), \mathbf{q}_h^k(\widehat{u}_h^k, f))$ by restricting the system (3.3) to an element.

Then using the superposition principle, the solution can be further split into two parts, namely

$$(u_h^k(\widehat{u}_h^k, f), \mathbf{q}_h^k(\widehat{u}_h^k, f)) = (u_h^k(\widehat{u}_h^k, 0), \mathbf{q}_h^k(\widehat{u}_h^k, 0)) + (u_h^k(0, f), \mathbf{q}_h^k(0, f)).$$

Then the Eq. (3.3c) reduces to finding $\widehat{u}_h^k \in M_h$ such that

$$a_h^{k-1}(\widehat{u}_h^k, \mu) = b_h^{k-1}(\mu), \quad \text{for all } \mu \in M_h.$$

where the bilinear form $a_h^{k-1}(\widehat{u}_h^k, \mu) : M_h \times M_h \rightarrow \mathbb{R}$ and the linear form $b_h^{k-1}(\mu) : M_h \rightarrow \mathbb{R}$ are defined as

$$a_h^{k-1}(\widehat{u}_h^k, \mu) = \langle \widehat{\mathbf{q}}_h^k(\widehat{u}_h^k, 0) \cdot \mathbf{n}, \mu \rangle_{\partial\mathcal{T}_h} \quad \text{and} \quad b_h^{k-1}(\mu) = -\langle \widehat{\mathbf{q}}_h^k(0, f) \cdot \mathbf{n}, \mu \rangle_{\partial\mathcal{T}_h},$$

respectively. By updating the coefficient value like $\alpha(\eta_h) = \alpha(u_h^k)$ and using iterative method, we can find $\widehat{u}_h^{k+1} \in M_h$ such that

$$a_h^k(\widehat{u}_h^{k+1}, \mu) = b_h^k(\mu), \quad \text{for all } \mu \in M_h.$$

For deriving a matrix equation, we insert (3.3e) and (3.3d) into (3.3a)-(3.3c) and obtain, by using integration by parts, that $(u_h^k, \mathbf{q}_h^k, \widehat{u}_h^k) \in W_h \times \mathbf{V}_h \times M_h$ is the solution of the following weak formulation:

$$a^{k-1}(\mathbf{q}_h^k, \mathbf{v}) - b(u_h^k, \mathbf{v}) + c(\widehat{u}_h^k, \mathbf{v}) = 0, \quad (4.1a)$$

$$-b(w, \mathbf{q}_h^k) - d(u_h^k, w) + e(\widehat{u}_h^k, w) = -f(w), \quad (4.1b)$$

$$c(\mu, \mathbf{q}_h^k) + e(\mu, u_h^k) - g(\mu, \widehat{u}_h^k) = 0, \quad (4.1c)$$

for all $(w, \mathbf{v}, \mu) \in W_h \times \mathbf{V}_h \times M_h$, $k \geq 1$, and the bilinear forms and the linear functional are defined by

$$\begin{aligned} a^{k-1}(\mathbf{q}, \mathbf{v}) &= (\alpha(u_h^{k-1})\mathbf{q}, \mathbf{v})_{\mathcal{T}_h}, & b(u, \mathbf{v}) &= (u, \nabla \cdot \mathbf{v})_{\mathcal{T}_h}, \\ c(\widehat{u}, \mathbf{v}) &= \langle \widehat{u}, \mathbf{v} \cdot \mathbf{n} \rangle_{\partial \mathcal{T}_h}, & d(u, w) &= \langle w, \tau u \rangle_{\partial \mathcal{T}_h}, \\ e(\mu, u) &= \langle \mu, \tau u \rangle_{\partial \mathcal{T}_h}, & g(\mu, \widehat{u}) &= \langle \mu, \tau \widehat{u} \rangle_{\partial \mathcal{T}_h}, & f(w) &= (f, w)_{\mathcal{T}_h}, \end{aligned} \quad (4.2)$$

for all $(u, \mathbf{q}, \widehat{u})$ and (w, \mathbf{v}, μ) in $W_h \times \mathbf{V}_h \times M_h$.

The discretization of the system of Eqs. (4.1) give rise to a matrix equation. The process of converting to a matrix equation follows the technique introduced in [26, 27]. Then, we have the following matrix equation:

$$\begin{bmatrix} A^{k-1} & -B^T & C^T \\ -B & -D & E \\ C & E^T & G \end{bmatrix} \begin{bmatrix} Q^k \\ U^k \\ \widehat{U}^k \end{bmatrix} = - \begin{bmatrix} 0 \\ F \\ 0 \end{bmatrix}. \quad (4.3)$$

Here Q^k, U^k , and \widehat{U}^k are the vectors of degrees of freedom for \mathbf{q}_h^k, u_h^k , and \widehat{u}_h^k , respectively. The matrices in (4.3) corresponding to the bilinear forms in (4.2) are in the order that they appear in the system (3.3).

Since the HDG method produces a final system in terms of globally coupled degrees of freedom of the numerical trace \widehat{u}_h^k (or \widehat{U}^k) only, the Eqs. (3.3a) and (3.3b) can be used to eliminate both \mathbf{q}_h^k and u_h^k in an element by element sense like (4.1). Then, we obtain a reduced globally coupled matrix equation only for \widehat{U}^k as

$$\mathbb{K}^{k-1} \widehat{U}^k = \mathbb{F}^{k-1} \quad (4.4)$$

where

$$\mathbb{K}^{k-1} = - \begin{bmatrix} C & E^T \end{bmatrix} \begin{bmatrix} A^{k-1} & -B^T \\ -B & -D \end{bmatrix}^{-1} \begin{bmatrix} C^T \\ E \end{bmatrix} + G,$$

and

$$\mathbb{F}^{k-1} = \begin{bmatrix} C & E^T \end{bmatrix} \begin{bmatrix} A^{k-1} & -B^T \\ -B & -D \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ F \end{bmatrix}.$$

We note that the matrix \mathbb{K}^{k-1} and the vector \mathbb{F}^{k-1} are associated with the bilinear form $a_h^{k-1}(\cdot, \cdot)$ and $b_h^{k-1}(\cdot)$, respectively.

After solving the Eq. (4.4), Q^k and U^k can be obtained from the following matrix equation.

$$\begin{bmatrix} Q^k \\ U^k \end{bmatrix} = \begin{bmatrix} A^{k-1} & -B^T \\ -B & -D \end{bmatrix}^{-1} \left(\begin{bmatrix} 0 \\ -F \end{bmatrix} - \begin{bmatrix} C^T \\ E \end{bmatrix} \widehat{U} \right).$$

4.2. Numerical examples. In this section, we present representative numerical examples in 2D to verify the optimal order of convergence for elliptic problems with nonlinear coefficients. For three-dimensional examples, we can identify them in the same way, but we omit them due to the complexity of numerical experiments. We are mainly interested in confirming the convergent ratio with various nonlinear coefficients that satisfy Assumptions 3.1. For each

nonlinear coefficient, we study the behavior of errors between the HDG solution and the exact solution when a mesh size changes.

We consider the domain $\Omega = (0, 1)^2$ and divide Ω into uniform triangulation consisting $2 \times 1/h \times 1/h$. We take the stabilization parameter $\tau = 1$ for the existence and uniqueness of the given system (3.3). The nonlinear coefficients for example 1, 2, and 3 are $\kappa_1(u) = \exp(u)$, $\kappa_2(u) = \exp(u^2)$, and $\kappa_3(u) = \frac{1}{1+u^2}$, respectively. We can easily check that these nonlinear terms satisfy Assumption 3.1 in the given domain. The source term f for all examples are chosen so that the exact solution is $u = \sin(\pi x)\sin(\pi y)$ in 2D. Figure 1 show the exact solution and nonlinear coefficients of examples.

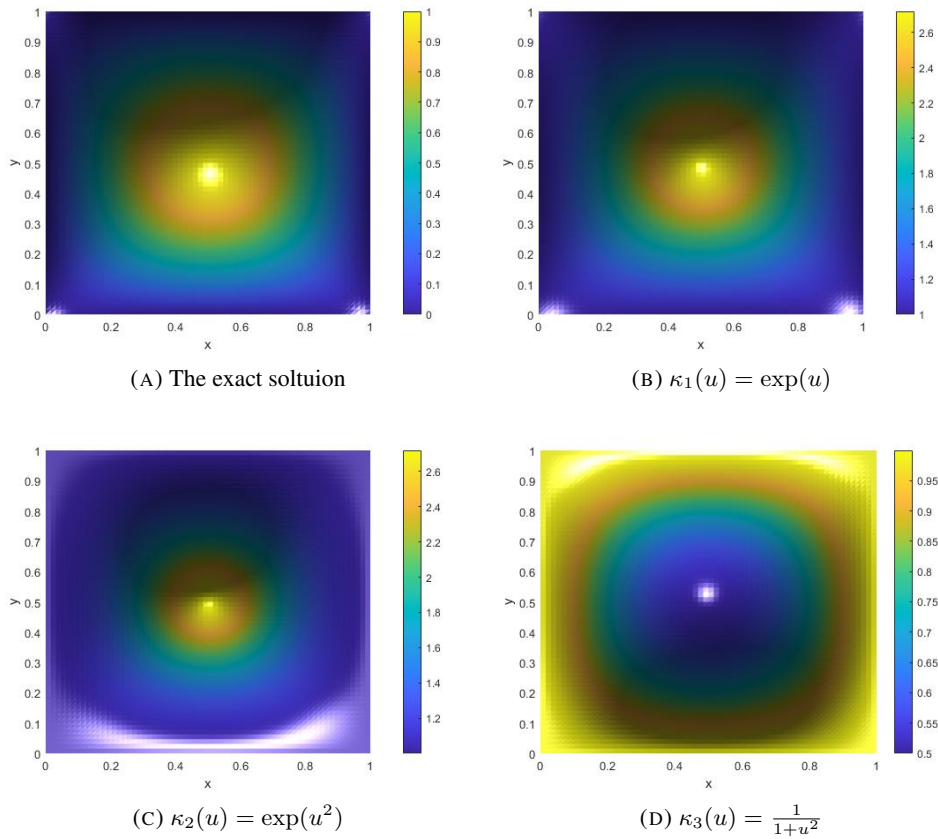


FIGURE 1. The exact solution and nonlinear coefficients for the example 1, 2, and 3.

We consider the following finite element spaces to apply the HDG method. W_h and V_h consist of piece-wise linear functions on \mathcal{T}_h , and M_h consists of piece-wise linear functions in

\mathcal{E}_h . Denote the exact solution and the HDG solution by (u^*, \mathbf{q}^*) and (u_h, \mathbf{q}_h) , respectively. We compute the relative L^2 error for the solution $\|u^* - u_h\|_{L^2(\Omega)}$ and for the flux $\|\mathbf{q}^* - \mathbf{q}_h\|_{L^2(\Omega)}$.

To solve the Eq. (4.4), we use iterative technique with given initial guess $\kappa^{-1}(u^0) = \alpha(u^0)$. In our experiment, we take the initial guess $\alpha(u^0) = 1$ and the tolerance $\delta = 10^{-8}$. The number of iterations for all examples presented in the paper is either 12 or 13.

h	$\ u^* - u_h\ _{L^2(\Omega)}$	order	$\ \mathbf{q}^* - \mathbf{q}_h\ _{L^2(\Omega)}$	order
1/4	0.4075	-	0.1987	-
1/8	0.1127	1.8543	0.0490	2.0197
1/16	0.0294	1.9386	0.0125	1.9709
1/32	0.0075	1.9709	0.0032	1.9658
1/64	0.0019	1.9809	0.0008	1.9894

TABLE 1. Errors and convergence orders for example 1

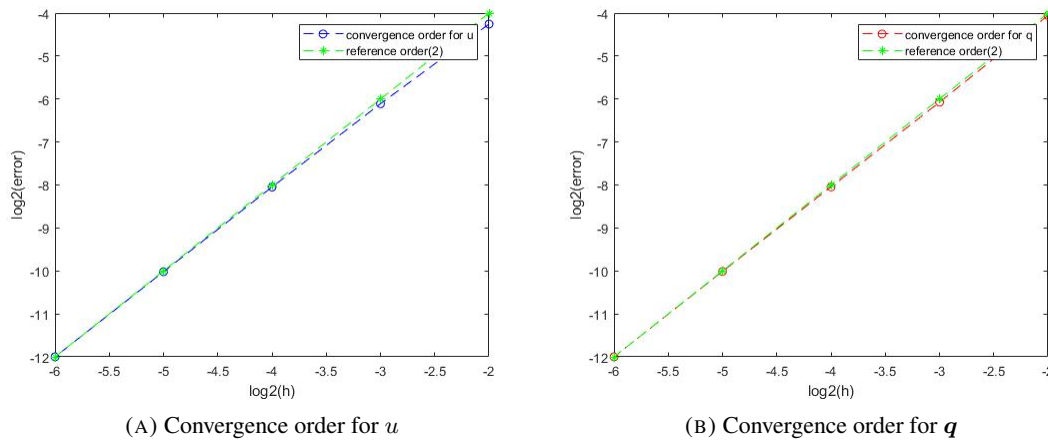


FIGURE 2. Log-log plots for example 1

Table 1 shows that the convergence order for u and \mathbf{q} is almost 2. This is consistent with the main results presented in the error analysis because we use piece-wise linear functions in the finite element spaces W_h , \mathbf{V}_h , and M_h . Also, Fig. 2 shows that the slope of two lines are almost the same.

For example 2, we consider the nonlinear coefficient $\kappa_2 = \exp(u^2)$. Table 2 and Fig. 3 show that the convergence behavior is similar to that of example 1.

For the last example, we take the nonlinear coefficient $\kappa_3 = \frac{1}{1+u^2}$. We observe similar accuracy in Table 3 and Fig. 4. This shows that the error analysis is correct.

h	$\ u^* - u_h\ _{L^2(\Omega)}$	order	$\ q^* - q_h\ _{L^2(\Omega)}$	order
1/4	0.3487	-	0.2065	-
1/8	0.0909	1.9396	0.0556	1.8930
1/16	0.0244	1.8974	0.0150	1.8901
1/32	0.0063	1.9535	0.0039	1.9434
1/64	0.0016	1.9773	0.0010	1.9793

TABLE 2. Errors and convergence orders for example 2

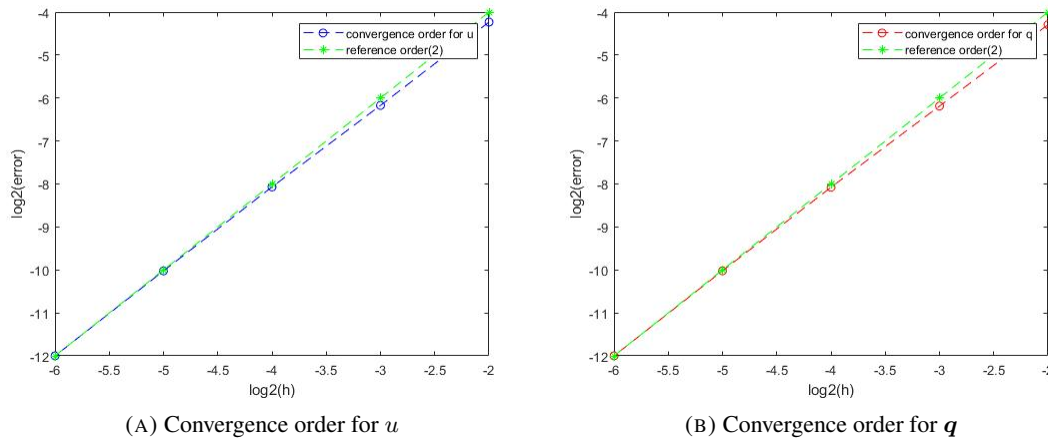


FIGURE 3. Log-log plots for example 2

h	$\ u^* - u_h\ _{L^2(\Omega)}$	order	$\ q^* - q_h\ _{L^2(\Omega)}$	order
1/4	0.1877	-	0.1510	-
1/8	0.0576	1.7043	0.0358	2.0765
1/16	0.0144	2.0000	0.0087	2.0409
1/32	0.0036	2.0000	0.0022	1.9835
1/64	0.0009	1.9910	0.0005	2.0339

TABLE 3. Errors and convergence orders for example 3

5. CONCLUSION

In this research, we analyze the HDG method for approximating the solution of elliptic PDEs with nonlinear coefficients. We assume that the nonlinear terms satisfy the Lipschitz and H^2 -regularity conditions to derive error estimates. Based on the projection analysis of the HDG method, we derive the optimal convergence ratio for the mesh size h . We also present matrix

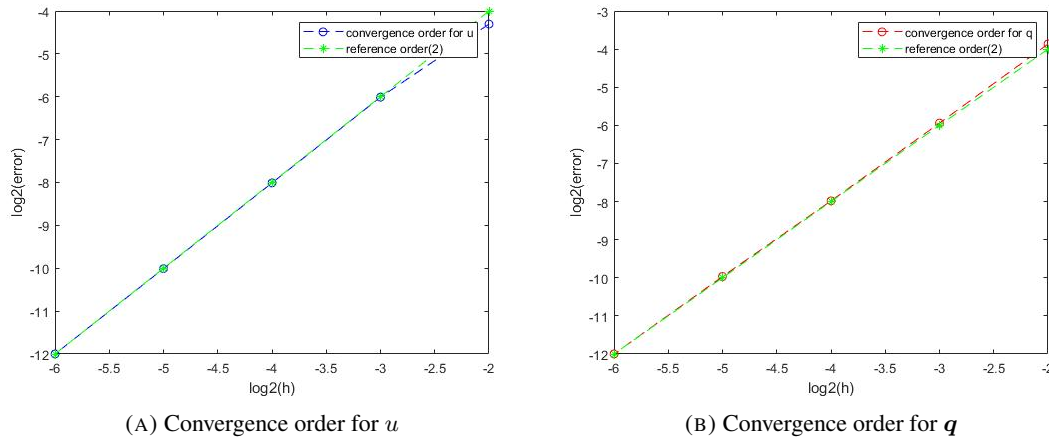


FIGURE 4. Log-log plots for example 3

formulations which can be easily applicable in the iteration procedure. We get the optimal order of convergence in the HDG method with nonlinear coefficients with suitable assumptions. We get the convergent ratio of $k + 1$ when we use the polynomials of degree k in finite element spaces for the HDG method. We present representative numerical examples that guarantee mathematical analysis. The results show the reliability and accuracy of error analysis. We plan to analyze the HDG and multiscale HDG method for many types of nonlinear PDEs using these approaches. These are practically applicable in many situations as flows in porous media.

ACKNOWLEDGEMENTS

This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (No. NRF-2020R1F1A1A01072414, NRF-2022R1F1A1069217).

REFERENCES

- [1] M. Moon. *Generalized multiscale hybridizable discontinuous Galerkin (GMsHDG) method for flows in nonlinear porous media*, Journal of Computational and Applied Mathematics, **415** (2022), 114440.
- [2] B. Cockburn, J. Gopalakrishnan, and R. Lazarov. *Unified hybridization of discontinuous Galerkin, mixed, and continuous Galerkin methods for second order elliptic problems*, SIAM Journal of Numerical Analysis, **47** (2009), 1319-1365.
- [3] D. N. Arnold, F. Brezzi, B. Cockburn, and L. D. Marini. *Unified analysis of discontinuous Galerkin methods for elliptic problems*, SIAM Journal of Numerical Analysis, **39** (2002), 1749-1779.
- [4] R. Kirby, S. Sherwin, and B. Cockburn. *To CG or to HDG : A comparative study*, Journal of Scientific Computing, **51** (2011), 183-212.
- [5] A. Huerta, A. Angeloski, X. Roca, and J. Peraire. *Efficiency of high-order elements for continuous and discontinuous Galerkin methods*, International Journal for Numerical Methods in Engineering, **96** (2013), 529-560.
- [6] G. Giorgiani, S. Fernandez-Mendez, and A. Huerta. *Hybridizable discontinuous Galerkin p -adaptivity for wave propagation problems*, International Journal for Numerical Methods in Fluids, **72** (2013), 1244-1262.

- [7] G. Giorgiani, S. Fernandez-Mendez, and A. Huerta. *Hybridizable discontinuous Galerkin with degree adaptivity for the incompressible Navier-Stokes equations*, Computers & Fluids, **98** (2014), 196-208.
- [8] S. Yakovlev, D. Moxey, R. M. Kirby, and S. J. Sherwin. *To CG or to HDG : A comparative study in 3D*, Journal of Scientific Computing, **67** (2016), 192-220.
- [9] B. Cockburn, W. Qiu, and K. Shi. *Conditions for superconvergence of HDG methods for second-order elliptic problems*, Mathematics of Computation, **81** (2012), 1327-1353.
- [10] B. Cockburn, J. Gopalakrishnan, and F.-J. Sayas, *A projection-based error analysis of HDG methods*, Mathematics of Computation, **79** (2010), 1351-1367.
- [11] N.C. Nguyen, J. Peraire, and B. Cockburn. *An implicit high-order hybridizable discontinuous Galerkin method for linear convection–diffusion equations*, Journal of Computational Physics, **228** (2009), 3232-3254.
- [12] N.C. Nguyen, J. Peraire, and B. Cockburn. *An implicit high-order hybridizable discontinuous Galerkin method for nonlinear convection–diffusion equations*, Journal of Computational Physics, **228** (2009), 8841-8855.
- [13] M. Moon, H.K. Jun, and T. Suh, *Error estimates on hybridizable discontinuous Galerkin methods for parabolic equations with nonlinear coefficients*, Advances in Mathematical Physics, **17** (2017), 1-11.
- [14] M. Stanglmeier, N.C. Nguyen, J. Peraire, and B. Cockburn. *An explicit hybridizable discontinuous Galerkin method for the acoustic wave equation*, Computer Methods in Applied Mechanics and Engineering, **300** (2016), 748-769.
- [15] M. Kronbichler, S. Schoeder, C. Müller, and W.A. Wall. *Comparison of implicit and explicit hybridizable discontinuous Galerkin methods for the acoustic wave equation*, International Journal for Numerical Methods in Engineering, **106** (2016), 712-739.
- [16] N.C. Nguyen, J. Peraire, and B. Cockburn. *A hybridizable discontinuous Galerkin method for Stokes flow*, Computer Methods in Applied Mechanics and Engineering, **199** (2010), 582-597.
- [17] G.N. Gatica and F.A. Sequeira. *Analysis of an augmented HDG method for a class of quasi-Newtonian Stokes flows*, Journal of Scientific Computing, **65** (2015), 1270-1308.
- [18] J. Peraire, N.C. Nguyen, and B. Cockburn. *A hybridizable discontinuous Galerkin method for the compressible Euler and Navier–Stokes equations*, AIAA, 48th AIAA Aerospace Sciences Meeting Including the New Horizons Forum and Aerospace Exposition, Florida 2010.
- [19] N.C. Nguyen, J. Peraire, and B. Cockburn. *A hybridizable discontinuous Galerkin method for the incompressible Navier-Stokes equations*, 48th AIAA Aerospace Sciences Meeting Including the New Horizons Forum and Aerospace Exposition. 2010.
- [20] N.C. Nguyen, J. Peraire, and B. Cockburn. *An implicit high-order hybridizable discontinuous Galerkin method for the incompressible Navier–Stokes equations*, Journal of Computational Physics, **230** (2011), 1147-1170.
- [21] A. Cesmelioglu, B. Cockburn, and W. Qiu. *Analysis of a hybridizable discontinuous Galerkin method for the steady-state incompressible Navier–Stokes equations*, Mathematics of Computation, **86** (2017), 1643-1670.
- [22] E.-J. Park. *Mixed finite element methods for nonlinear second-order elliptic problems*, SIAM Journal on Numerical Analysis, **32** (1995), 865-885.
- [23] D. Kim and E.-J. Park. *A priori and a posteriori analysis of mixed finite element methods for nonlinear elliptic equations*, SIAM Journal on Numerical Analysis, **48** (2010), 1186-1207.
- [24] Y. Sangita, A. Pani, and E.-J. Park. *Superconvergent discontinuous Galerkin methods for nonlinear elliptic equations*, Mathematics of Computation, **82** (2013), 1297-1335.
- [25] A. Muhammad, E.-J. Park, and D. Shin. *Analysis of multiscale mortar mixed approximation of nonlinear elliptic equations*, Computers and Mathematics with Applications, **75** (2018), 401-418.
- [26] R. Sevilla and A. Huerta. *Tutorial on Hybridizable Discontinuous Galerkin (HDG) for second-order elliptic problems*, Advanced finite element technologies 566, Springer, Cham, 2016.
- [27] M. Moon and Y. H. Lim. *Superconvergence of Hybridizable Discontinuous Galerkin method for second-order elliptic equations*, Journal of the Korean Society for Industrial and Applied Mathematics, **20** (2016), 295-308.
- [28] P. Grisvard. *Elliptic problems in nonsmooth domains*, Classics in Applied Mathematics 69, Pitman, Boston, MA, 1985.

- [29] P. Knabner. *Numerical methods for elliptic and parabolic partial differential equations*, Springer, New York, 2003.
- [30] L.C. Evans. *Partial differential equations*, Graduate Studies on Mathematics 19, American Mathematical Society, 2010.