# PERFORMANCE OF DISCONTINUOUS GALERKIN METHODS FOR ELLIPTIC PDES*

PAUL CASTILLO†

**Abstract.** In this paper, we compare the performance of several discontinuous Galerkin (DG) methods for elliptic partial differential equations (PDEs) on a model problem. Theoretical estimates of the condition number of the stiffness matrix are given for DG methods whose bilinear form is symmetric and which are shown to be numerically sharp. Then the efficiency of the methods in the computation of both the potential and its gradient is tested on unstructured triangular meshes.

**Key words.** finite element, discontinuous Galerkin methods, conditioning of stiffness matrix

**AMS subject classification.** 65N30

**PII.** S1064827501388339

**1. Introduction.** In the last decade, several discontinuous Galerkin (DG) methods have been proposed for solving nonlinear hyperbolic and convection dominated problems; see [7] for an introduction to the subject and [14] for an overview of the state of the art. DG methods are preferred over traditional finite volume methods because they provide high-order accurate approximations; they have a high degree of parallelism; and, since no interelement continuity is imposed, polynomials of arbitrary degree can be used on different elements, making these methods suitable for *hp* refinement. Over the past years, there has been a tremendous interest in their application to problems in which the diffusion is not negligible and to pure elliptic problems.

Recently, Arnold et al. [3] developed a unified framework in which theoretical stability analysis and optimal error estimates can be obtained for virtually all existing DG methods. However, they do not discuss important issues that could be relevant to the practitioner. In this paper, we complete the work presented in [3] by analyzing the methods from a practical point of view. The DG methods we consider are the following: The Babuška–Zlámal penalty method [5], which is the simplest of all DG methods; the *interior penalty* (IP) method [16, 6, 24, 2], which is one of the first symmetric DG methods for linear and nonlinear parabolic problems with provable optimal error estimates; the steady state version of the so-called *local discontinuous Galerkin* (LDG) method for purely elliptic problems [15, 11, 12]; and, finally, a class of nonsymmetric methods called *nonsymmetric interior penalty Galerkin* (NIPG) [21, 22, 23] which include the method proposed by Baumann and Oden [9].

We are interested in the quality and efficiency of the numerical approximation. We compare the above methods with respect to asymptotic behavior of the spectral condition number of the stiffness matrix, storage cost, rates of convergence, and accuracy of the approximation of the potential and gradient. This comparison is carried out on a model elliptic problem with a smooth solution. The spectral condition number is analyzed numerically and theoretically for symmetric DG methods. The

†Center for Applied Scientific Computing, Lawrence Livermore National Laboratory, 7000 East Avenue, P.O. Box 808 L-551, Livermore, CA 94551 (castillo17@llnl.gov).

analysis aims to get explicit expressions of the bounds in terms of the stabilization parameters of each method. For the nonsymmetric methods, we perform a numerical study of the spectral condition number as a function of the mesh size as well as of their stabilization parameters.

The organization of the paper is as follows. In section 2, we present the general formulation of a class of DG methods which contains all the methods considered in this paper. In section 3, we present a theoretical analysis of the spectral condition number for DG methods with symmetric bilinear forms. We also carry out a numerical study of the condition number in terms of the stabilization parameter of each method. In section 4, we perform a comparison from a practical perspective. We analyze the storage cost and compare the accuracy of the potential and the gradients; unstructured meshes are used. Finally, we end in section 5 with some concluding remarks.

**2. General formulation of DG methods.** We describe the formulation of a general DG method applied to the elliptic model problem with Dirichlet boundary conditions

$$(2.1) \qquad -\Delta u = f \quad \text{in } \Omega,$$

$$(2.2) \qquad u = g \quad \text{on } \partial\Omega,$$

where $\Omega$ is a bounded convex domain in $\mathbb{R}^d$.

Based on [11, 3], we introduce a new variable $\boldsymbol{q} = \nabla u$. We can rewrite our model problem as a system of the form

$$(2.3) \qquad \boldsymbol{q} = \nabla u \quad \text{in } \Omega,$$

$$(2.4) \qquad -\nabla \cdot \boldsymbol{q} = f \quad \text{in } \Omega,$$

$$u = g \quad \text{on } \partial\Omega.$$

Let $\mathcal{T}_h$ be a general triangulation of $\Omega$. The weak formulation is obtained by multiplying (2.3) and (2.4) by smooth test functions $\boldsymbol{r}$ and $v$, respectively, on each element $T$ of $\mathcal{T}_h$. After integrating by parts, we obtain the weak formulation

$$(2.5) \qquad \int_T \boldsymbol{q} \cdot \boldsymbol{r} = \oint_{\partial T} u\boldsymbol{r} \cdot \vec{n}_T - \int_T u\nabla \cdot \boldsymbol{r},$$

$$(2.6) \qquad \int_T \boldsymbol{q} \cdot \nabla v = \oint_{\partial T} v\boldsymbol{q} \cdot \vec{n}_T + \int_T fv,$$

where $\vec{n}_T$ is the outward unit vector normal to the element $T$. Note that the above equations are well defined for any functions $(u, \boldsymbol{q})$ and $(v, \boldsymbol{r})$ in $\mathcal{V} \times \mathcal{M}$, where

$$\mathcal{V} = \{u \in L_2(\Omega) : u\big|_T \in H^1(T) \ \forall T \in \mathcal{T}_h\},$$
$$\mathcal{M} = \{\boldsymbol{q} \in (L_2(\Omega))^d : \boldsymbol{q}\big|_T \in H^1(T)^d \ \forall T \in \mathcal{T}_h\}.$$

Next, we seek to approximate the exact solution $(u, \boldsymbol{q})$ with functions $(u_h, \boldsymbol{q}_h)$ in the finite element space $\mathcal{V}_{\mathbf{h}} \times \mathcal{M}_{\mathbf{h}} \subset \mathcal{V} \times \mathcal{M}$, where

$$\mathcal{V}_{\mathbf{h}} = \{u \in L_2(\Omega) : u\big|_T \in \mathcal{P}_k(T) \ \forall T \in \mathcal{T}_h\},$$
$$\mathcal{M}_{\mathbf{h}} = \{\boldsymbol{q} \in (L_2(\Omega))^d : \boldsymbol{q}\big|_T \in \mathcal{P}_k(T)^d \ \forall T \in \mathcal{T}_h\},$$

and the *local* finite element space $\mathcal{P}_k(T)$ is the set of polynomials of degree at most $k$. The finite element solution $(u_h, \boldsymbol{q}_h)$ is defined by using the aforementioned weak

formulation by requiring that for all $T \in \mathcal{T}_h$, and for all $(v, \boldsymbol{r}) \in \mathcal{P}_k(T) \times \mathcal{P}_k(T)^d$, we have

$$(2.7) \qquad \int_T \boldsymbol{q_h} \cdot \boldsymbol{r} = \oint_{\partial T} \widehat{u_h}_{\{e,T\}} \boldsymbol{r} \cdot \vec{n}_T - \int_T u_h \nabla \cdot \boldsymbol{r},$$

$$(2.8) \qquad \int_T \boldsymbol{q_h} \cdot \nabla v = \oint_{\partial T} v \widehat{\boldsymbol{q_h}}_{\{e,T\}} \cdot \vec{n}_T + \int_T fv,$$

where $\widehat{u_h}_{\{e,T\}}$ and $\widehat{\boldsymbol{q_h}}_{\{e,T\}}$ are the so-called *numerical fluxes* which can be thought of as approximations to the traces of the function $u$ and $\boldsymbol{q}$, respectively. We assume that these fluxes are local quantities in the sense that they depend only on the traces to the edge $e$ of functions $u_{h|\{T,K\}}$, $q_{h|\{T,K\}}$, and/or $\nabla u_{h|\{T,K\}}$, where $T$ and $K$ are the elements sharing edge $e$. Moreover, we expect these functions to satisfy some basic properties such as *consistency*, that is, $\widehat{u_h}_{\{e,T\}} = u$ and $\widehat{\boldsymbol{q_h}}_{\{e,T\}} = \nabla u$, for a smooth function $u$, required in numerical methods for conservation laws, and the so-called *conservation* property, which can be formally defined as follows: Let $T$ and $K$ be the elements sharing edge $e$. A numerical flux $\widehat{\sigma}_{e,T}$ is conservative if

$$\widehat{\sigma}_{e,T} = \widehat{\sigma}_{e,K}.$$

By suitably choosing the numerical fluxes, we obtain the DG methods we are interested in, as shown in [3]. To define the fluxes, we need to introduce some notation. Let $e$ be an interior edge shared by elements $T$ and $K$; we denote by $\vec{n}_T$ and $\vec{n}_K$ the outward unit normal vectors on $e$, relative to $T$ and $K$, respectively. For any function $v \in \mathcal{V}$, we define the jump, $[\![v]\!]$, and the average, $\{\!\{v\}\!\}$, of $v$ on an interior edge $e$ by

$$[\![v]\!] = v_{|K} \vec{n}_K + v_{|T} \vec{n}_T \quad \text{and} \quad \{\!\{v\}\!\} = \frac{1}{2} \left( v_{|K} + v_{|T} \right).$$

For the boundary edges we simply define $[\![u]\!] = u_{|T} \vec{n}_T$ and $\{\!\{u\}\!\} = u_{|T}$. For any function $r \in \mathcal{M}$, the jump and the average are defined similarly:

$$[\![\boldsymbol{r}]\!] = \boldsymbol{r}_{|K} \cdot \vec{n}_K + \boldsymbol{r}_{|T} \cdot \vec{n}_T \quad \text{and} \quad \{\!\{\boldsymbol{r}\}\!\} = \frac{1}{2} \left( \boldsymbol{r}_{|K} + \boldsymbol{r}_{|T} \right).$$

In Table 2.1, we show the definition of the numerical fluxes for the DG methods considered in this paper. Observe that, in the Babuška–Zlámal method, $\widehat{u_h}_{\{e,T\}}$ is not conservative and $\widehat{\boldsymbol{q_h}}_{\{e,T\}}$ is not consistent. Although optimal error estimates still can be obtained by using penalty terms of the order $O(h^{-(2p+1)})$, this method is not suitable for practical computations, since the condition number of the stiffness matrix is proportional to $O(h^{-(2p+2)})$. In the IP method both numerical fluxes are consistent and conservative. The method is symmetric and achieves optimal rates of convergence for both the potential and the gradient using, unlike the Babuška–Zlámal method, a penalization term independent of the approximation polynomial degree. However, the stabilization parameter $\eta$ is mesh-dependent and must be chosen large enough to make the bilinear form coercive. Like the IP method, LDG is a symmetric method and both numerical fluxes are consistent and conservative. Unlike the IP method, the LDG method is stable for $\eta > 0$. In the NIPG and Baumann–Oden methods, the numerical flux $\widehat{u_h}_{\{e,T\}}$ is not conservative; this renders the bilinear form nonsymmetric. Note that the Baumann–Oden method does not require a penalty term. However, this lack of stabilization is responsible for the suboptimality of the accuracy of the method. This behavior has been observed numerically in the $L_2$ norm. The NIPG method of Rivière, Wheeler, and Girault [21, 22] tries to fix this problem by including a penalty term of order $O(h^{-\beta})$. In this paper we denote by NIPG1 and NIPG3 the NIPG methods with a penalty term of order $O(h^{-1})$ and $O(h^{-3})$, respectively.

TABLE 2.1

*Definition of the numerical fluxes for various DG methods, when using approximations of degree p.*

| Method | $\widehat{u_h}_{\{e,T\}}$ | $\widehat{\boldsymbol{q_h}}_{\{e,T\}}$ |
|---|---|---|
| Babuška–Zlámal | $u_{h|T}$ | $-\frac{\eta}{h_e^{2p+1}}[\![u_h]\!]$ |
| IP | $\{\!\{u_h\}\!\}$ | $\{\!\{\nabla u_h\}\!\} - \frac{\eta}{h_e}[\![u_h]\!]$ |
| LDG | $\{\!\{u_h\}\!\} + \beta_e \cdot [\![u_h]\!]$ | $\{\!\{\boldsymbol{q_h}\}\!\} - \beta_e[\![\boldsymbol{q_h}]\!] - \frac{\eta}{h_e}[\![u_h]\!]$ |
| Baumann–Oden | $\{\!\{u_h\}\!\} + \vec{n}_T \cdot [\![u_h]\!]$ | $\{\!\{\nabla u_h\}\!\}$ |
| NIPG1 | $\{\!\{u_h\}\!\} + \vec{n}_T \cdot [\![u_h]\!]$ | $\{\!\{\nabla u_h\}\!\} - \frac{\eta}{h_e}[\![u_h]\!]$ |
| NIPG3 | $\{\!\{u_h\}\!\} + \vec{n}_T \cdot [\![u_h]\!]$ | $\{\!\{\nabla u_h\}\!\} - \frac{\eta}{h_e^3}[\![u_h]\!]$ |

**3. Conditioning of the stiffness matrix.** It is well known that, for the standard finite element method and under certain conditions of the uniformity of the mesh, the spectral condition number of the stiffness matrix, i.e., the ratio between the largest and smallest eigenvalue, is of order $O(h^{-2})$, where $h$ is the mesh size. In this section, we analyze the spectral condition number of the reduced stiffness matrix, i.e., the matrix obtained after the elimination of the auxiliary variable $\boldsymbol{q_h}$. We derive theoretical bounds for the DG methods with symmetric bilinear forms. The proof relies on a Poincaré–Friedrichs-type inequality; see Lemma 2.2 of Arnold [2]. For nonsymmetric methods, namely, the method of Oden, Babuška, and Baumann [19], the NIPG method of Rivière, Wheeler, and Girault [21, 22], and the method of Süli, Schwab, and Houston [23], a theoretical characterization of the spectrum is significantly difficult; this is why a numerical study of the spectral condition number is mandatory.

**3.1. Stiffness matrix.** The global linear system obtained from the general discontinuous formulation is a large sparse system which involves not only the potential but also the auxiliary variable $\boldsymbol{q}$. Moreover, the resulting system is indefinite. For these reasons it is preferable in practice to eliminate the auxiliary variable and solve only for the potential. We must point out that this elimination is possible only because of the particular choice of the numerical flux $\widehat{u_h}_{\{e,T\}}$, which depends only on the variable $u$.

The assembly of the reduced stiffness matrix can be done either by computing a Schur complement from the global linear system or by exploiting the discontinuous formulation at the element level. The simplest approach is the classical procedure used in mixed methods. First, we obtain the matrix expression for (2.7) and (2.8),

$$(3.1) \qquad MQ_h + B_1 U_h = f,$$
$$(3.2) \qquad B_2 Q_h + C U_h = g,$$

where $M$ is the mass matrix, $B_1$ and $B_2$ are the gradient and divergence operators, $C$ is the stabilizing term, and $Q_h$ and $U_h$ are the discrete versions of the functions $\boldsymbol{q_h}$ and $u_h$, respectively. Since $M$ is block diagonal, we easily can get an explicit expression for $Q_h$ from (3.1). Using this expression in (3.2), we obtain the following

linear system for the variable $U_h$:

$$\left(C - B_2 M^{-1} B_1\right) U_h = g - B_2 M^{-1} f.$$

**3.2. Preliminaries.** For each element $T \in \mathcal{T}_h$, we denote by $h_T$ the diameter of $T$, by $\rho_T$ the diameter of the largest ball contained in $T$, and set $h = \max\{h_T : T \in \mathcal{T}_h\}$. We denote by $\mathcal{E}_i$ the set of interior edges and by $\mathcal{E}$ the set of all the edges, including the boundary edges. We assume that the triangulation is quasi-regular; that is, there exists a positive constant $\sigma$ such that

$$\forall\, T \in \mathcal{T}_h, \quad \frac{h_T}{\rho_T} \leq \sigma.$$

Let $\{\phi_i\}$ be a basis of $\mathcal{V}_\mathbf{h}$; we denote by $|u|_\Omega = \left(\sum_i \gamma_i^2\right)^{1/2}$ the $l_2$ norm in $\mathcal{V}_\mathbf{h}$, where $u = \sum_i \gamma_i \phi_i$. Then we have the following relation between the $l_2$ norm and the $L_2$ norm $\|\cdot\|_{0,\Omega}$ in $\mathcal{V}_\mathbf{h}$; see [20, Prop. 6.3.1]. There exist positive constants $C^*$ and $C^{**}$ that depend on $d$ and $\sigma$ such that for each $u \in \mathcal{V}_\mathbf{h}$ we have

$$C^* h^d |u|_\Omega^2 \leq \|u\|_{0,\Omega}^2 \leq h^d C^{**} |u|_\Omega^2.$$

Let $\mathcal{A}\left(\cdot,\cdot\right)$ be an arbitrary symmetric bilinear form, and let $A$ be its associated stiffness matrix. The spectral condition number $\kappa(A)$ can then be bounded by

$$(3.3) \qquad\qquad\qquad\qquad \kappa(A) \leq R^* \frac{C_2}{C_1},$$

where $R^*$ is the ratio $C^{**}/C^*$ and $C_1, C_2$ are positive constants, which can depend on the mesh size $h$, such that for all $u \in \mathcal{V}_\mathbf{h}$, $\mathcal{A}\left(\cdot,\cdot\right)$ is bounded by

$$C_1 \|u\|_{0,\Omega}^2 \leq \mathcal{A}\left(u,u\right) \leq C_2 \|u\|_{0,\Omega}^2.$$

Thus, to obtain the bound (3.3) we prove the existence of the constants $C_1$ and $C_2$. We recall the following standard inverse and trace inequalities.

LEMMA 3.1. *For any $u \in \mathcal{P}_k(T)$ there exist positive constants $C_1 = C_1(k,\sigma)$ and $C_2 = C_2(k,\sigma)$ such that*

$$(3.4) \qquad\qquad\qquad\qquad \|\nabla u\|_{0,T}^2 \leq C_1 h_T^{-2} \|u\|_{0,T}^2,$$
$$(3.5) \qquad\qquad\qquad\qquad \|u\|_{0,e}^2 \leq C_2 h_T^{-1} \|u\|_{0,T}^2,$$

*where $e$ is an edge of $\partial T$.*

As mentioned before, our results rely on a Poincaré–Friedrichs type of inequality. This inequality will be crucial in proving the coercivity of the bilinear form of several DG methods. Here, we want to slightly refine that result in order to obtain an explicit expression in terms of the parameters of a DG method.

LEMMA 3.2 (see Lemma 2.2 in [2]). *Let $u$ be a function in $\mathcal{V}$. For any positive number $\lambda$, there exist constants $C_1$ and $C_2$ such that*

$$\|u\|_{0,\Omega}^2 \leq C_\lambda \left(\sum_{T \in \mathcal{T}_h} \|\nabla u\|_{0,T}^2 + \int_\mathcal{E} \frac{\lambda}{h_e} [\![u]\!]^2\right),$$

*where $C_\lambda = C_1 + \frac{\max\{1,h^2\}}{\lambda} C_2$; $C_1$ depends on $\Omega$ and $C_2$ depends on the minimum angle bound of the mesh.*

*Proof.* We proceed as in [2]. Let $z \in H^2(\Omega) \bigcap H_0^1(\Omega)$ be the solution of the problem; $-\Delta z = u$ in $\Omega$ and $z = 0$ on $\partial\Omega$. It is clear that $(z, \nabla z)$ satisfies

$$\|u\|_{0,\Omega}^2 = \sum_{T \in \mathcal{T}_h} \int_T \nabla u \cdot \nabla z - \int_{\mathcal{E}} [\![u]\!] \nabla z.$$

Then after applying the Cauchy–Schwarz inequality we obtain

$$\|u\|_{0,\Omega}^2 \leq \left( \sum_{T \in \mathcal{T}_h} \|\nabla u\|_{0,T}^2 + \int_{\mathcal{E}} \frac{\lambda}{h_e} [\![u]\!]^2 \right)^{1/2} \left( \sum_{T \in \mathcal{T}_h} \|\nabla z\|_{0,T}^2 + \int_{\mathcal{E}} \frac{h_e}{\lambda} \left\| \frac{\partial z}{\partial n} \right\|_{0,e}^2 \right)^{1/2},$$

where $\lambda$ is an arbitrary positive parameter. From the regularity of $z$ [18] there exists a positive constant $C_1 = C_1(\Omega)$ such that

$$\sum_{T \in \mathcal{T}_h} \|\nabla z\|_{0,T}^2 \leq C_1 \|u\|_{0,\Omega}^2.$$

Moreover, using the trace inequality [1, 2]

$$\forall \phi \in H^2(T), \quad \left\| \frac{\partial \phi}{\partial n} \right\|_{0,e}^2 \leq C_2 \left( \frac{1}{h_e} |\phi|_{1,T}^2 + h_e |\phi|_{2,T}^2 \right),$$

where $|\phi|_{k,T}^2 = \|D^k \phi\|_{0,T}^2$ and the constant $C_2 = C_2(\sigma)$ depends only on the minimal angle bound, we get

$$\sum_{T \in \mathcal{T}_h} \|\nabla z\|_{0,T}^2 + \int_{\mathcal{E}} \frac{h_e}{\lambda} \left\| \frac{\partial z}{\partial n} \right\|_{0,e}^2 \leq C_\lambda \|u_h\|_{0,\Omega}^2,$$

where $C_\lambda = C_1 + \frac{\max\{1, h^2\}}{\lambda} C_2$, and where $C_1, C_2$ have the above stated dependence. This completes the proof. $\square$

We are now ready to prove the estimates of the condition number for those methods with symmetric bilinear form. In subsection 3.3, we prove that for the Babuška–Zlámal penalty method the condition number is of order $O(h^{-(2p+2)})$, where $p$ is the approximation polynomial degree, showing that the method is not suitable for high-order approximations. In subsections 3.4 and 3.5, we analyze the IP and LDG methods, respectively. We show that both methods have a condition number of order $O(h^{-2})$, which is similar to the asymptotic behavior of a standard continuous finite element. Finally, we end this section with a numerical study of the condition number for the nonsymmetric methods (subsection 3.6), where we show that the condition number is of order $O(h^{-2})$ for the NIPG1 and Baumann–Oden methods and of $O(h^{-4})$ for the NIPG3 method.

**3.3. The Babuška–Zlámal penalty method.** To obtain the bilinear form we proceed as follows. First, observe that the auxiliary variable $q_h$ is equal to $\nabla u_h$. Indeed, for any $u_h \in \mathcal{V}_\mathbf{h}$, using integration by parts in the right term of (2.7) and the definition of $\widehat{u_{h\{e,T\}}}$, we have

$$\forall r \in \mathcal{M}_h, \quad \int_T (q_h - \nabla u_h) \cdot r = \oint_{\partial T} \left( \widehat{u_{h\{e,T\}}} - u_h \right) r \cdot \vec{n}_T = 0.$$

Hence, $\boldsymbol{q_h} = \nabla u_h$. Inserting this expression into (2.8) and adding over all the elements, we obtain

$$\mathcal{A}\left(u_h, v_h\right) := \sum_{T \in \mathcal{T}_h} \int_T \nabla u_h \cdot \nabla v_h + \int_{\mathcal{E}} \frac{\eta}{h^{2p+1}} [\![u_h]\!][\![v_h]\!] = \int_{\Omega} f v_h.$$

The corresponding energy norm $\|\!|\cdot|\!\|$ is given by

(3.6) $$\|\!|u_h|\!\|^2 = \mathcal{A}\left(u_h, u_h\right) = \sum_{T \in \mathcal{T}_h} \|\nabla u_h\|_{0,T}^2 + \int_{\mathcal{E}} \frac{\eta}{h^{2p+1}} [\![u_h]\!]^2.$$

The following result shows the dependence of the condition number of $\mathcal{A}\left(\cdot, \cdot\right)$ with the mesh size $h$.

THEOREM 3.3. *For $h < 1$, we have*

$$C_{\eta,h}^1 \|u_h\|_{0,\Omega}^2 \leq \mathcal{A}\left(u_h, u_h\right) \leq C_{\eta,h}^2 \|u_h\|_{0,\Omega}^2,$$

*where $C_{\eta,h}^1 = 1/(C_1^* + \frac{h^{2p}}{\eta} C_2^*)$ and $C_{\eta,h}^2 = \frac{1}{h^2} C_1 + \frac{\eta}{h^{2p+2}} C_2$ in which $C_1^*, C_2^*, C_1, C_2$ are positive constants that depend on $\Omega, \sigma, k$.*

*Proof.* The lower bound constant $C_{\eta,h}^1$ can be obtained by using the Poincaré–Friedrichs inequality (Lemma 3.2) with $\lambda = \frac{\eta}{h^{2p}}$. We have

$$\|u\|_{0,\Omega}^2 \leq C_\lambda \left( \sum_{T \in \mathcal{T}_h} \|\nabla u\|_{0,T}^2 + \int_{\mathcal{E}} \frac{\eta}{h^{2p+1}} [\![u]\!]^2 \right),$$

where $C_\lambda = C_1^* + \frac{h^{2p}}{\eta} C_2^*$, $C_1^* = C(\Omega)$, and $C_2^* = C(\sigma)$. Then set $C_{\eta,h}^1 = 1/(C_1^* + \frac{h^{2p}}{\eta} C_2^*)$.

The upper bound constant $C_{\eta,h}^2$ can be obtained easily after an application of inequalities (3.4) and (3.5):

$$\mathcal{A}\left(u_h, u_h\right) \leq \left( \frac{C_1}{h^2} + \frac{\eta C_2}{h^{2p+2}} \right) \|u_h\|_{0,\Omega}^2,$$

where $C_1 = C(k, \sigma)$ and $C_2 = C(k, \sigma)$. Hence, let $C_{\eta,h}^2 = C_1/h^2 + \eta C_2/h^{2p+2}$. □

Using (3.3), it is then easy to see that if $A_h$ is the corresponding stiffness matrix, then the spectral condition number $\kappa(A_h)$ is bounded by

$$\kappa(A_h) \leq R^* \frac{C_{\eta,h}^2}{C_{\eta,h}^1} = R^* \left( \frac{1}{h^2} C_1 + \frac{\eta}{h^{2p+2}} C_2 \right) \left( C_1^* + \frac{h^{2p}}{\eta} C_2^* \right) = O\left(h^{-2p-2}\right).$$

**3.4. The IP method.** We now consider the IP method. It can be shown [2] that the energy norm $\|\!|\cdot|\!\|$ associated with this method is defined as

$$\|\!|u_h|\!\|^2 = \mathcal{A}_h\left(u_h, u_h\right) = \sum_{T \in \mathcal{T}_h} \|\nabla u_h\|_{0,T}^2 - 2\int_{\mathcal{E}} [\![u_h]\!]\{\!\{\nabla u_h\}\!\} + \int_{\mathcal{E}} \frac{\eta}{h_e} [\![u_h]\!]^2.$$

THEOREM 3.4. *We have*

$$C_\eta^1 \|u_h\|_{0,\Omega}^2 \leq \mathcal{A}\left(u_h, u_h\right) \leq \frac{C_\eta^2}{h^2} \|u_h\|_{0,\Omega}^2,$$

where $C_\eta^1 = 1/(C_1^* + \frac{C_2^*}{\eta - C_3})$, $C_\eta^2 = C_1 + C_2\eta$, and where $C_1, C_2, C_3$ are generic constants that depend on $\Omega, \sigma, k$.

This represents a major drawback for the IP method since the stabilization parameter $\eta$ must be larger than the constant $C_3$, which depends on the shape regularity of the mesh and the approximation polynomial degree and is difficult to find in practice.

*Proof.* Using the inequality $ab \le \frac{\epsilon}{2}a^2 + \frac{1}{2\epsilon}b^2$, for any $\epsilon > 0$ we have

$$\int_{\mathcal{E}} [\![u_h]\!] \{\!\!\{\nabla u_h\}\!\!\} \le \frac{\epsilon}{2} \int_{\mathcal{E}} \frac{1}{h} [\![u_h]\!]^2 + \frac{1}{2\epsilon} \int_{\mathcal{E}} h \, \|\{\!\!\{\nabla u_h\}\!\!\}\|_{0,e}^2 \,,$$

and by the trace inequality (3.5) we obtain

(3.7) $$\int_{\mathcal{E}} [\![u_h]\!] \{\!\!\{\nabla u_h\}\!\!\} \le \frac{\epsilon}{2} \int_{\mathcal{E}} \frac{1}{h} [\![u_h]\!]^2 + \frac{C}{2\epsilon} \sum_{T \in \mathcal{T}_h} \|\nabla u_h\|_{0,T}^2 \,,$$

where $C = C(k, \sigma)$ is a positive constant. Then

$$\mathcal{A}(u_h, u_h) \ge (1 - C/\epsilon) \sum_{T \in \mathcal{T}_h} \|\nabla u_h\|_{0,T}^2 + \int_{\mathcal{E}} \left(\frac{\eta - \epsilon}{h}\right) [\![u_h]\!]^2,$$

$$\ge (1 - C/\epsilon) \left(\sum_{T \in \mathcal{T}_h} \|\nabla u_h\|_{0,T}^2 + \frac{\eta - \epsilon}{1 - C/\epsilon} \int_{\mathcal{E}} \frac{1}{h} [\![u_h]\!]^2\right).$$

The existence of the constant $C_\eta^1$ will follow by choosing $\eta$ such that $C < \epsilon < \eta$. The lower bound is obtained by using Lemma 3.2 with $\lambda = (\eta - \epsilon)/(1 - C/\epsilon)$:

$$\mathcal{A}(u_h, u_h) \ge \frac{1 - C/\epsilon}{C_\lambda} \|u_h\|_{0,\Omega}^2 \,.$$

Then we have $C_\eta^1 = \frac{1 - C/\epsilon}{C_\lambda}$, which is of the form $1/(C_1^* + \frac{C_2^*}{\eta - \epsilon})$.

To obtain the upper bound, we proceed as follows:

$$\mathcal{A}_h(u_h, u_h) \le (1 + C) \sum_{T \in \mathcal{T}_h} \|\nabla u_h\|_{0,T}^2 + (1 + \eta) \int_{\mathcal{E}} \frac{1}{h} [\![u_h]\!]^2 \quad \text{by (3.7) with } \epsilon = 1,$$

$$\le \frac{C_1}{h^2} \|u_h\|_{0,\Omega}^2 + \frac{C_2(1 + \eta)}{h^2} \|u_h\|_{0,\Omega}^2 \quad \text{by inequalities (3.4) and (3.5)},$$

$$\le \frac{C_\eta^2}{h^2} \|u_h\|_{0,\Omega}^2 \,,$$

where $C_\eta^2 = C_1 + \eta C_2$ and where $C_1 = C(k, \sigma)$, $C_2 = C(k, \sigma)$ are positive constants. $\square$

An upper bound for the spectral condition number of the stiffness matrix of the IP method is given by

(3.8) $$\kappa(A_h) \le R^* (C_1 + \eta C_2) \left(C_1^* + \frac{C_2^*}{\eta - C}\right) \frac{1}{h^2},$$

where $C, C_1, C_2, C_1^*, C_2^*$ are positive constants.

In Figure 1, we show the condition number of the reduced stiffness matrix for the IP method as a function of mesh size $h$. A given line represents the variation of the conditioning for a fixed polynomial degree. Using linear regression we find that
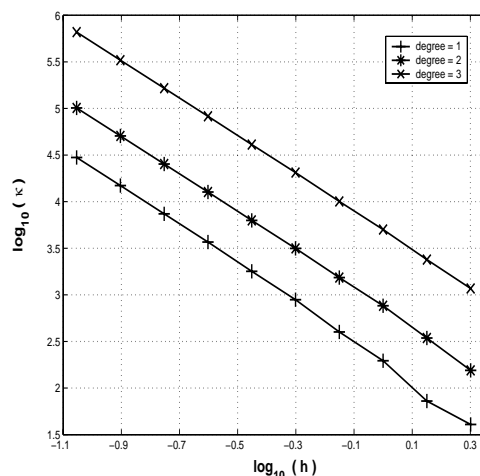
FIG. 1. *Spectral condition number of the IP method as a function of mesh size h for the model problem* (2.1) *in* $\Omega = (-1,1) \times (-1,1)$.

the slopes of the lines are $-2.1360$ for $p = 1$, $-2.0603$ for $p = 2$, and $-2.0276$ for $p = 3$. Thus, the conditioning of the matrix behaves like $O(h^{-2})$ independently of the polynomial degree, showing that our analysis is sharp.

The upper bound (3.8) shows that there exists a mesh-dependent positive constant $\eta_0$ such that the spectral condition number $\kappa(\eta)$ has the following asymptotic behavior:

$$\kappa(\eta) = \begin{cases} O\left(1/(\eta - \eta_0)\right) & \text{if } \eta - \eta_0 \ll 1, \\ O(\eta) & \text{if } \eta \gg 1. \end{cases}$$

In Figure 2, we show the function $\kappa(\eta)$ for linear and quadratic approximations. The plot is in logarithmic scale and has been shifted to the critical values $\eta_0 \approx 2.8666825$ for $p = 1$ and $\eta_0 \approx 6.9016989$ for $p = 2$. Observe that $\kappa(\eta)$ grows linearly as $\eta$ goes to infinity and when $\eta$ approaches $\eta_0$, that is, $\log_{10}(\eta - \eta_0) \longrightarrow -\infty$; the curve resembles a straight line with slope $-1$, meaning that $\kappa(\eta) = O\left(1/(\eta - \eta_0)\right)$. This shows that our theoretical analysis is also sharp for the stabilization parameter $\eta$.

**3.5. The LDG method.** The LDG method was introduced by Cockburn and Shu in [15] as a generalization of the DG method proposed by Bassi and Rebay [8] for the solution of the compressible Navier–Stokes equations. Recently, Castillo et al. [11] presented the first optimal a priori error estimates for the steady state version of the LDG method on arbitrary meshes. Using the general formulation of a DG method from (2.7) and (2.8), it can be shown [11] that a general expression for the mesh-dependent energy norm $\|\cdot\|$ is given by

$$(3.9) \qquad \|u_h\|^2 = \mathcal{A}_h\left(u_h, u_h\right) = \sum_{T \in \mathcal{T}_h} \|\boldsymbol{q_h}\|_{0,T}^2 + \int_{\mathcal{E}} \frac{\eta}{h_e} [\![u_h]\!]^2,$$

where the vector function $\boldsymbol{q_h}$ belongs to $\mathcal{M}_{\mathbf{h}}$ and is defined by 2.7 in terms of $u_h$. In general, for any DG method defined by (2.7) and (2.8), we can view $\boldsymbol{q_h}$ as an approximation of the gradient of the exact solution $u$. This approximation is of practical interest since it could lead to a more accurate approximation of the gradient or could even converge with rates higher than that of the standard piecewise gradient
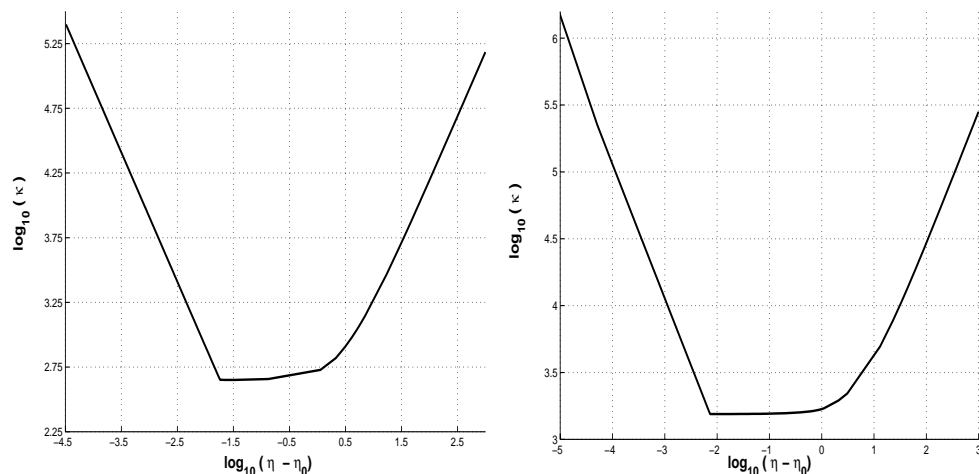
FIG. 2. *Spectral condition number of the IP method, using linear (left) and quadratic (right) approximations as a function of the stabilization parameter $\eta$, for the model problem (2.1) in $\Omega = (-1, 1) \times (-1, 1)$.*

$\nabla u_h$. Recently, Cockburn et al. [13] obtained a superconvergence result for the gradient on Cartesian grids, using $\boldsymbol{q_h}$ as an approximation to $\nabla u$. The following lemma establishes a connection between $\boldsymbol{q_h}$ and $\nabla u_h$.

LEMMA 3.5. *For any function $u_h \in \mathcal{V}_\mathbf{h}$, there exists a positive constant $C_\beta = C(\beta, k, \sigma)$ such that*

$$\sum_{T \in \mathcal{T}_h} \|\boldsymbol{q_h} - \nabla u_h\|_{0,T}^2 \leq C_\beta \int_{\mathcal{E}_i} \frac{1}{h_e} [\![u_h]\!]^2.$$

*Proof.* Since $\boldsymbol{q_h}$ and $\nabla u_h$ belong to $\mathcal{M}_\mathbf{h}$, we can take $\boldsymbol{r} = \boldsymbol{q_h} - \nabla u_h$ in (2.7). Then after integrating by parts in the right term of (2.7) and adding over all the elements, we obtain

$$\sum_{T \in \mathcal{T}_h} \|\boldsymbol{q_h} - \nabla u_h\|_{0,T}^2 = \int_{\mathcal{E}_i} [\![(\widehat{u_{h}}_{\{e,T\}} - u_h)(\boldsymbol{q_h} - \nabla u_h)]\!]$$
$$= \int_{\mathcal{E}_i} (\beta_e [\![\boldsymbol{q_h} - \nabla u_h]\!] - \{\!\{\boldsymbol{q_h} - \nabla u_h\}\!\}) \cdot [\![u_h]\!],$$

since

$$[\![u\boldsymbol{r}]\!] = \{\!\{u\}\!\}[\![\boldsymbol{r}]\!] + [\![u]\!]\{\!\{\boldsymbol{r}\}\!\}.$$

Using the Cauchy–Schwarz inequality, we get

$$\sum_{T \in \mathcal{T}_h} \|\boldsymbol{q_h} - \nabla u_h\|_{0,T}^2 \leq \left( \int_{\mathcal{E}_i} h_e \, \|\beta_e [\![\boldsymbol{q_h} - \nabla u_h]\!] - \{\!\{\boldsymbol{q_h} - \nabla u_h\}\!\}\|_{0,e}^2 \right)^{1/2} \left( \int_{\mathcal{E}_i} \frac{1}{h_e} [\![u_h]\!]^2 \right)^{1/2}.$$

Finally, after applying the inverse inequality (3.5), we obtain the estimate

$$\sum_{T \in \mathcal{T}_h} \|\boldsymbol{q_h} - \nabla u_h\|_{0,T}^2 \leq C_\beta \int_{\mathcal{E}_i} \frac{1}{h_e} [\![u_h]\!]^2,$$

where the constant $C_\beta$ has the stated dependence. This completes the proof. $\square$

The following result shows that the spectral condition number of the stiffness matrix for the LDG method is bounded by $O(h^{-2})$.

THEOREM 3.6. *We have*

$$C^1_{\beta,\eta} \|u_h\|^2_{0,\Omega} \le \mathcal{A}_h(u_h, u_h) \le \frac{C^2_{\beta,\eta}}{h^2} \|u_h\|^2_{0,\Omega},$$

*where $C^1_{\beta,\eta} = 1/(C_\beta \max\{1, \frac{1}{\eta}\})$, $C^2_{\beta,\eta} = C_1 + \eta C_2$, and where $C_\beta, C_1, C_2$ are positive constants that depend on $\beta, \sigma, k$.*

*Proof.* To obtain the upper bound, we proceed as follows:

$$\mathcal{A}_h(u_h, u_h) \le 2 \sum_{T \in \mathcal{T}_h} \|\nabla u_h\|^2_{0,T} + 2 \sum_{T \in \mathcal{T}_h} \|\boldsymbol{q_h} - \nabla u_h\|^2_{0,T} + \int_{\mathcal{E}} \frac{\eta}{h_e} [\![u_h]\!]^2$$

$$\le 2 \sum_{T \in \mathcal{T}_h} \|\nabla u_h\|^2_{0,T} + C_\beta \int_{\mathcal{E}} \frac{1}{h_e} [\![u_h]\!]^2 + \int_{\mathcal{E}} \frac{\eta}{h_e} [\![u_h]\!]^2 \quad \text{by Lemma 3.5}$$

$$\le \frac{C_1}{h^2} \|u_h\|^2_{0,\Omega} + (C_\beta + \eta) \frac{C_2}{h^2} \|u_h\|^2_{0,\Omega} \quad \text{by inequalities (3.4) and (3.5)}$$

$$\le \frac{C^2_{\beta,\eta}}{h^2} \|u_h\|^2_{0,\Omega},$$

where the positive constant $C^2_{\beta,\eta}$ has the general form $C^2_{\beta,\eta} = C_1 + \eta C_2$ and the positive generic constants $C_1$ and $C_2$ have dependence $C_1 = C(k, \sigma, \beta)$ and $C_2 = C(k, \sigma)$.

To prove the lower bound we proceed as follows:

$$\|u\|^2_{0,\Omega} \le C \left( \|\nabla u_h\|^2_{0,\Omega} + \int_{\mathcal{E}} \frac{1}{h_e} [\![u]\!]^2 \right) \quad \text{by Lemma 3.2 with } \lambda = 1$$

$$\le C \left( 2 \sum_{T \in \mathcal{T}_h} \|\boldsymbol{q_h}\|^2_{0,T} + 2 \sum_{T \in \mathcal{T}_h} \|\nabla u_h - \boldsymbol{q_h}\|^2_{0,T} + \int_{\mathcal{E}} \frac{1}{h_e} [\![u]\!]^2 \right),$$

$$\le C \left( 2 \sum_{T \in \mathcal{T}_h} \|\boldsymbol{q_h}\|^2_{0,T} + (1 + C_\beta) \int_{\mathcal{E}} \frac{1}{h_e} [\![u]\!]^2 \right) \quad \text{by Lemma 3.5}$$

$$\le C_2 \left( \sum_{T \in \mathcal{T}_h} \|\boldsymbol{q_h}\|^2_{0,T} + \int_{\mathcal{E}} \frac{1}{h_e} [\![u]\!]^2 \right),$$

where the positive constant $C_2 = \max\{2, 1 + C_\beta\}$. Then we have

$$\|u\|^2_{0,\Omega} \le C_2 \max\left\{1, \frac{1}{\eta}\right\} \left( \sum_{T \in \mathcal{T}_h} \|\boldsymbol{q_h}\|^2_{0,T} + \int_{\mathcal{E}} \frac{\eta}{h} [\![u_h]\!]^2 \right).$$

The lower bound constant follows by setting $C^1_{\beta,\eta} = 1/(C_\beta \max\{1, \frac{1}{\eta}\})$. $\square$

The spectral condition number $\kappa(A_h)$ of the stiffness matrix $A_h$ of the LDG method is bounded by

$$(3.10) \qquad \kappa(A_h) \le R^* \frac{C^2_{\beta,\eta}}{C^1_{\beta,\eta}} h^{-2} = R^* (C_1 + \eta C_2) \left( C_\beta \max\left\{1, \frac{1}{\eta}\right\} \right) h^{-2}.$$
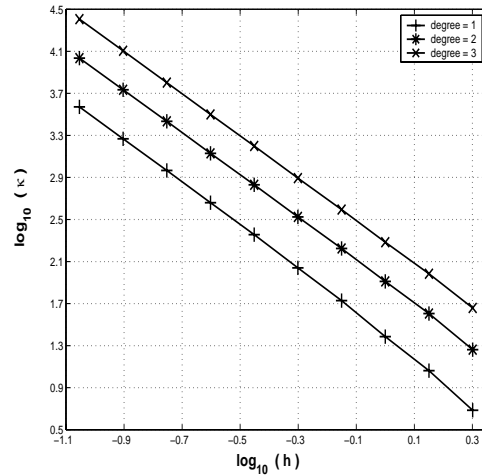
FIG. 3. *Spectral condition number of the LDG method as a function of mesh size h for the model problem* (2.1) *in* $\Omega = (-1, 1) \times (-1, 1)$.
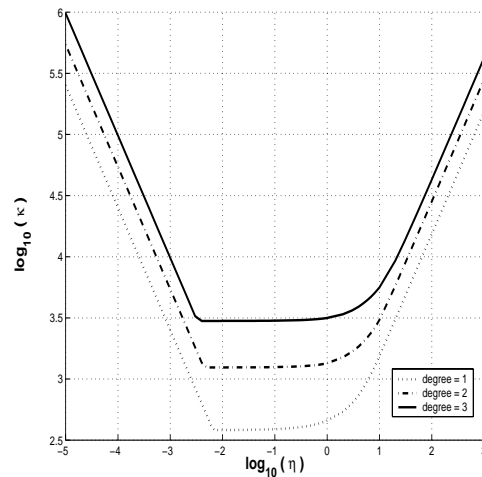


FIG. 4. *Spectral condition number of the LDG method for approximations of degree* $p = 1$, 2, *and* 3 *for the model problem* (2.1) *in* $\Omega = (-1, 1) \times (-1, 1)$.

In Figure 3, we show the condition number of the reduced stiffness matrix for the LDG method as a function of mesh size $h$. Each line represents the condition number using a different polynomial degree $p = 1, 2, 3$. Using linear regression we find that the slopes of the lines are $-2.1109$ for $p = 1$, $-2.0334$ for $p = 2$, and $-2.0207$ for $p = 3$; thus the condition number varies as $O(h^{-2})$ independently of the polynomial degree, as predicted by our theoretical analysis.

Note that, unlike the IP method, the LDG method is stable for any $\eta > 0$. Next we show that the theoretical upper bound obtained in (3.10) is sharp. In Figure 4, we plot the spectral condition number $\kappa(\eta)$ as a function of $\eta$ on a logarithmic scale; we have used approximations of degree $p = 1, 2, 3$. Observe that the global asymptotic behavior is independent of the approximation polynomial degree. The stiffness matrix

becomes ill-conditioned as $\eta$ approaches 0, meaning that the LDG method becomes less stable, as expected. For $\eta \ll 1$, $\log \kappa(\eta)$ behaves like a straight line with slope $-1$, that is, $\kappa(\eta) = O(\frac{1}{\eta})$. Finally, for $\eta \gg 1$, the curve behaves like a straight line with slope 1, that is, $\kappa(\eta) = O(\eta)$. Both asymptotic behaviors agree with our theoretical analysis.

**3.6. The NIPG methods.** This class of DG methods includes the methods analyzed by Rivière, Wheeler, and Girault [21, 22], the Baumann–Oden method [4, 9], and the discontinuous $hp$-finite element method of Süli, Schwab, and Houston [23].

Unlike the IP method, the Baumann–Oden method is always well defined and weakly stable in the sense used in [3]. Unfortunately, this lack of stabilization has a negative impact on the accuracy of the method since, as has been observed, a loss in the order of convergence for polynomials of even degree occurs on meshes with quadrilateral elements [4]. In our experiments, using unstructured triangular meshes, we also have observed numerically this loss of accuracy in the $L_2$ norm [10]. The NIPG methods of Rivière, Wheeler, and Girault [21, 22] are stabilizations of the Baumann–Oden method. Here we consider two cases: NIPG1 (i.e., $\beta = 1$, which corresponds to a stabilization term similar to that of the LDG and IP methods) and NIPG3 (i.e., $\beta = 3$, which is the minimal value that gives optimal rates of convergence in the $L_2$ norm for the potential. The analysis can be found in [21, 22]).

In Table 3.1, we present the spectral condition number of the stiffness matrix for the Baumann–Oden, NIPG1, and NIPG3 methods on a sequence of nested structured meshes. Since the Baumann–Oden method is stable for polynomials of degree greater than or equal to 2, the column corresponding to $p = 1$ is left empty. We assume that the asymptotic behavior of the spectral condition number is of order $O(h^\alpha)$. In Table 3.2, we show a numerical estimate of $\alpha$, which was obtained by linear regression. The spectral condition number for the NIPG1 and Baumann–Oden methods is of order $O(h^{-2})$ and of order $O(h^{-4})$ for the NIPG3 method.

TABLE 3.1

*Spectral condition number for a sequence of nested structured meshes for the model problem (2.1) in $\Omega = (-1, 1) \times (-1, 1)$.*

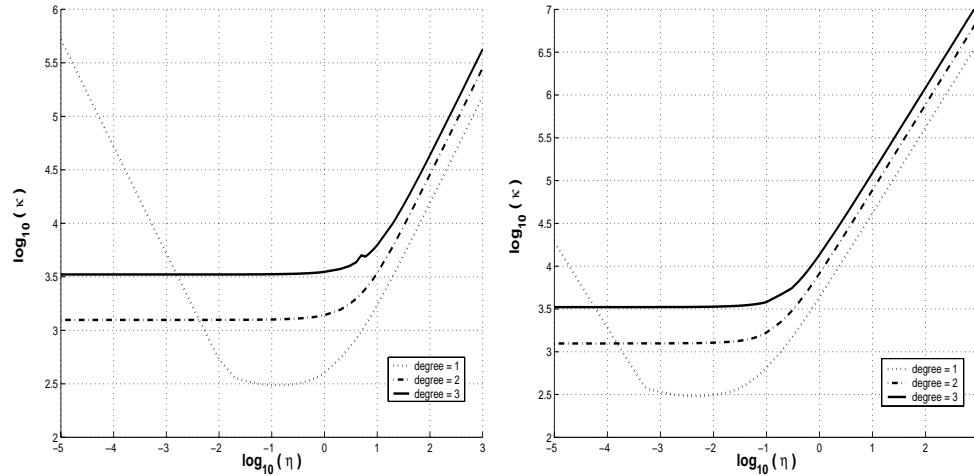|              | $p = 1$    | $p = 2$    | $p = 3$    |
| ------------ | ---------- | ---------- | ---------- |
| Baumann–Oden | —          | 7.8607e+01 | 2.0620e+02 |
|              | —          | 3.1155e+02 | 8.2943e+02 |
|              | —          | 1.2509e+03 | 3.3246e+03 |
|              | —          | 5.0113e+03 | 1.3318e+04 |
|              | —          | 2.0076e+04 | 5.3348e+04 |
| NIPG1        | 8.0255e+02 | 1.7005e+03 | 2.6404e+03 |
|              | 3.7107e+03 | 7.0816e+03 | 1.0785e+04 |
|              | 1.5485e+04 | 2.8607e+04 | 4.3358e+04 |
|              | 6.2639e+04 | 1.1471e+05 | 1.7365e+05 |
|              | 2.5338e+05 | 4.6000e+05 | 6.9546e+05 |
| NIPG3        | 1.3979e+03 | 2.9992e+03 | 4.8334e+03 |
|              | 2.5128e+04 | 4.8332e+04 | 7.6398e+04 |
|              | 4.1695e+05 | 7.7360e+05 | 1.2185e+06 |
|              | 6.7348e+06 | 1.2375e+07 | 1.9481e+07 |
|              | 1.0878e+08 | 1.9796e+08 | 3.1146e+08 |

FIG. 5. *Spectral condition number as a function of the stabilization parameter $\eta$ for the NIPG method, NIPG1 $\beta = 1$ (left), and NIPG3 $\beta = 3$ (right) for the model problem* (2.1) *in* $\Omega = (-1, 1) \times (-1, 1)$.

TABLE 3.2
*Numerical estimate of order $\alpha$ of the conditioning number for the sequence of nested meshes.*

|               | $p = 1$       | $p = 2$       | $p = 3$       |
|---------------|---------------|---------------|---------------|
| Baumann–Oden  | —             | -2.0001e+00   | -2.0036e+00   |
| NIPG1         | -2.0682e+00   | -2.0177e+00   | -2.0091e+00   |
| NIPG3         | -4.0562e+00   | -4.0021e+00   | -3.9946e+00   |

In Figure 5, we plot $\kappa(\eta)$ for the NIPG1 (left) and NIPG3 (right) methods for linear, quadratic, and cubic approximations. Since the NIPG method is a stabilization of the Baumann–Oden method, the spectral condition number should remain bounded as $\eta$ goes to 0 for polynomials of degree greater than or equal to 2. For linear approximations, the Baumann–Oden method is unstable; thus the condition number of the NIPG method should grow as $\eta$ becomes smaller. Now, for large values of $\eta$, $\kappa$ grows linearly with respect to $\eta$ for both values of $\beta$. This suggests an asymptotic behavior of order $O(\eta)$ for $\eta \gg 1$, which is similar to that of the IP and LDG methods.

**4. Comparison of methods.** In this section we compare storage cost, condition number, rates of convergence, and accuracy. To this end, we solve the model problem (2.1) with homogeneous boundary conditions in the convex domain, $\Omega = (-1, 1) \times (-1, 1)$. The right-hand side is chosen such that the exact solution is the smooth function $u(x, y)$ given by

$$u(x, y) = 4(1 - x^2)(1 - y^2)e^{0.75(x+y)}.$$

We use a restarted GMRES method to solve nonsymmetric linear systems and use the conjugate gradient method for the others. In order to obtain as much precision as possible, the stopping criterion is such that the relative residual norm is less than $10^{-13}$.

**4.1. Memory requirements.** The size of the stencil provides an upper bound for the storage cost of the stiffness matrix and has a direct impact on parallel implementations. It is completely determined by the choice of the numerical flux $\widehat{\boldsymbol{q_h}}_{\{e,T\}}$. If this flux depends on the auxiliary variable $\boldsymbol{q_h}$, then an element interacts with its immediate neighbors and also with the neighbors of its neighbors; thus the maximum size of the stencil is 10 (in two-dimensional domains). This is the case for the LDG method. Otherwise, an element interacts only with its immediate neighbors, and hence the possible maximum size for a stencil is 4, i.e., the IP, Baumann–Oden, and NIPG methods.

We measure the storage cost of the stiffness matrix in terms of $nnz$, the total number of nonzero entries. Let $s$ be the size of the stencil, $d_k$ the dimension of the local polynomial space (we assume that the polynomial degree $k$ is the same on each element), and $n$ the total number of elements in the triangulation $\mathcal{T}_h$. Then an upper bound for $nnz$ is

$$nnz \leq nsd_k^2.$$

However, we must point out that this is just a rough estimate. In practice the actual number is below this bound and depends on the polynomial basis. Here we have used the high-order orthogonal basis for triangular elements proposed by Dubiner [17]. In Table 4.1 we show the normalized storage cost. For any given mesh and polynomial degree, we compute the total number of nonzero entries of the stiffness matrix of each method and normalize it against the minimum storage. Since the IP, NIPG, and Baumann–Oden methods have a compact stencil, this ratio should be close to 1. However, for the LDG method the ratio should be close to the theoretical bound $(10/4) = 2.5$.

TABLE 4.1
*Normalized storage cost.*

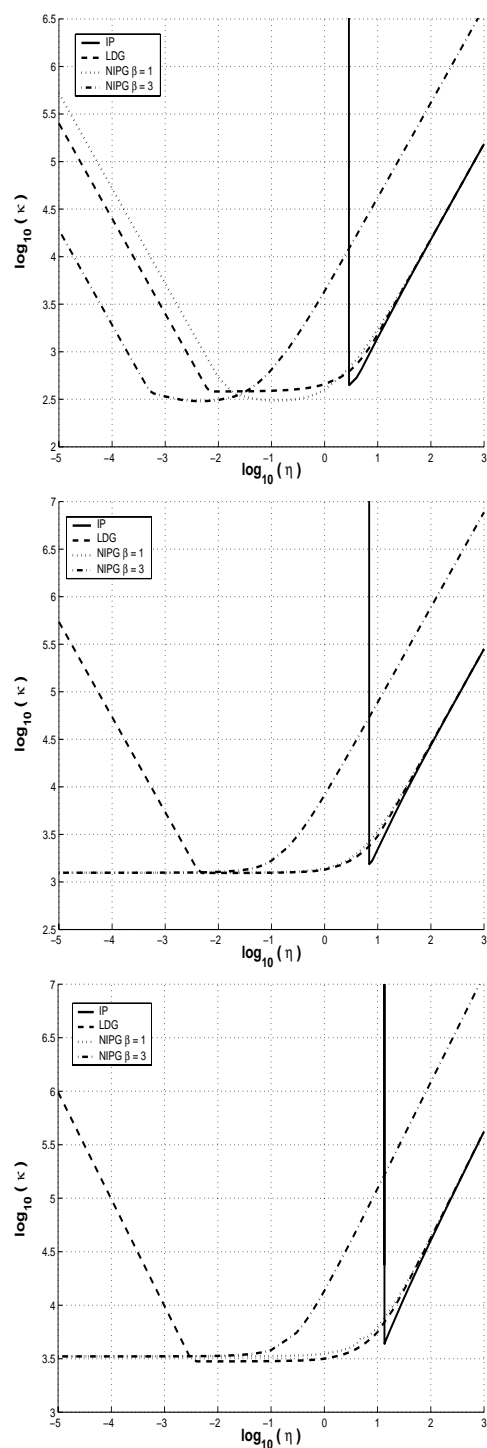|  | dofs | IP | LDG | BO | NIPG1 | NIPG3 |
|---|---|---|---|---|---|---|
| | $16 \times 3$ | 1.000 | 1.571 | —— | 1.000 | 1.041 |
| | $64 \times 3$ | 1.000 | 1.861 | —— | 1.000 | 1.059 |
| $p = 1$ | $256 \times 3$ | 1.000 | 2.007 | —— | 1.000 | 1.068 |
| | $1024 \times 3$ | 1.000 | 2.081 | —— | 1.000 | 1.073 |
| | $4096 \times 3$ | 1.000 | 2.117 | —— | 1.000 | 1.075 |
| | $16 \times 6$ | 1.183 | 1.917 | 1.000 | 1.195 | 1.207 |
| | $64 \times 6$ | 1.185 | 2.288 | 1.000 | 1.202 | 1.219 |
| $p = 2$ | $256 \times 6$ | 1.186 | 2.465 | 1.000 | 1.206 | 1.225 |
| | $1024 \times 6$ | 1.186 | 2.552 | 1.000 | 1.207 | 1.228 |
| | $4096 \times 6$ | 1.187 | 2.595 | 1.000 | 1.208 | 1.229 |
| | $16 \times 10$ | 1.090 | 1.792 | 1.000 | 1.106 | 1.110 |
| | $64 \times 10$ | 1.087 | 2.159 | 1.000 | 1.110 | 1.116 |
| $p = 3$ | $256 \times 10$ | 1.086 | 2.336 | 1.000 | 1.112 | 1.119 |
| | $1024 \times 10$ | 1.085 | 2.423 | 1.000 | 1.113 | 1.120 |
| | $4096 \times 10$ | 1.085 | 2.466 | 1.000 | 1.114 | 1.121 |

FIG. 6. *Comparison of the spectral condition number $\kappa(\eta)$ as a function of the stabilization parameter $\eta$ for linear (top), quadratic (middle), and cubic (bottom) approximations for the model problem (2.1) in $\Omega = (-1, 1) \times (-1, 1)$.*
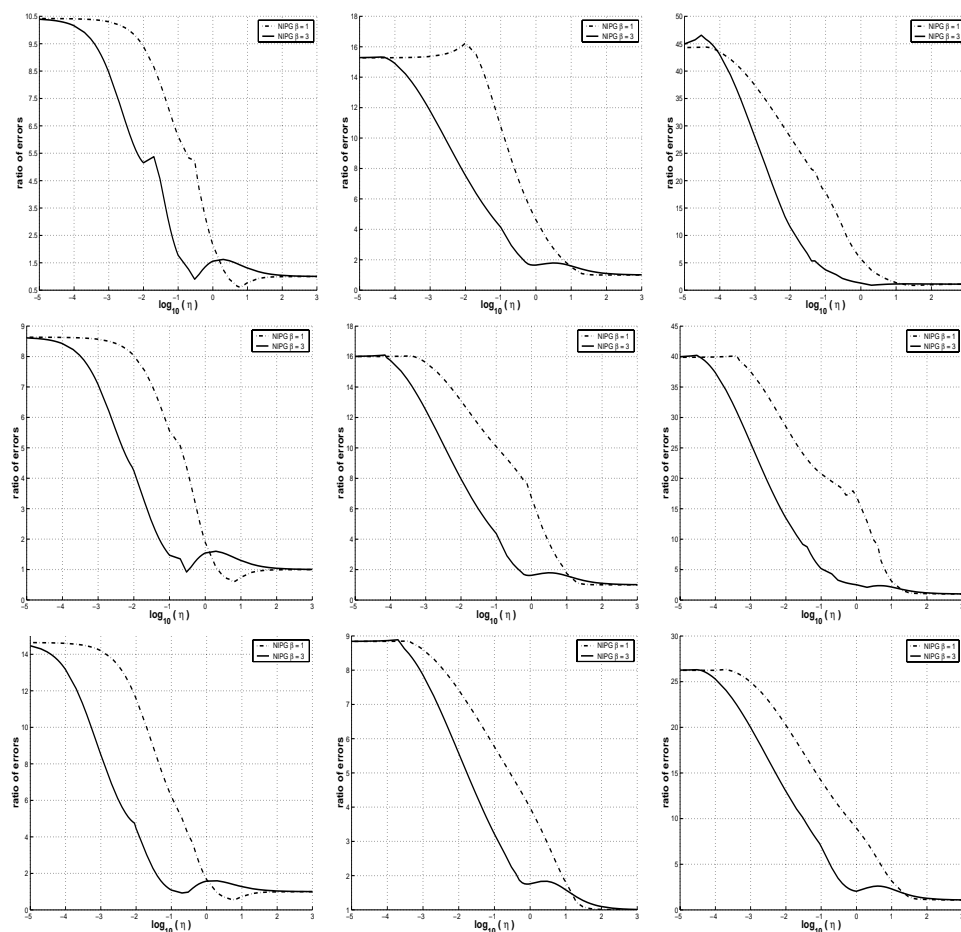
FIG. 7. *Ratio of the errors for the potential NIPG : LDG. Row i corresponds to mesh i and column j to approximation of degree j.*

**4.2. Conditioning.** In Table 4.2, we summarize the theoretical estimates obtained in section 3. We show the order of the condition number as well as the order of the penalization term of the method as a function of mesh size $h$ and polynomial degree $p$. We also include the nonsymmetric methods based on our numerical results. For small mesh size $h$, the penalization term dominates; thus for a penalization term of order $O(h^\alpha)$, the spectral condition number should be of order $O(h^{\alpha+1})$.

We consider the condition number as a function of the stabilization parameter $\eta$. In Figure 6, we plotted the spectral condition number of the IP, LDG, NIPG1, and NIPG3 methods for linear, quadratic, and cubic approximations on a structured triangular mesh with 256 elements. For large values of $\eta$, the condition numbers of the IP, LDG, and NIPG1 methods are asymptotically the same. This result is, to some extent, expected since their stabilization terms are of the same order, i.e., $O(\frac{\eta}{h})$. This suggests that when no preconditioning is used, the IP method will take less CPU time than the LDG method, since performing a matrix-vector multiplication with the LDG method is 2 to 2.5 times more expensive than with the IP. For the nonsymmetric methods, the NIPG1 method will outperform the NIPG3 method, since the condition
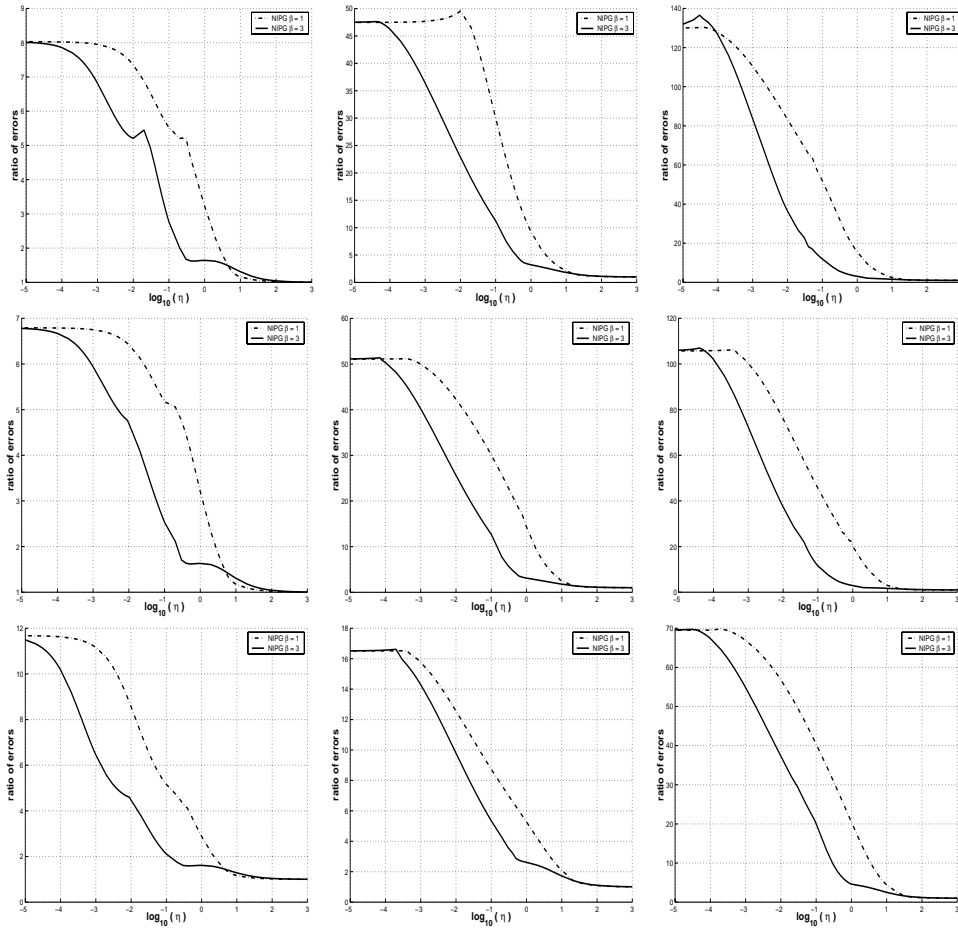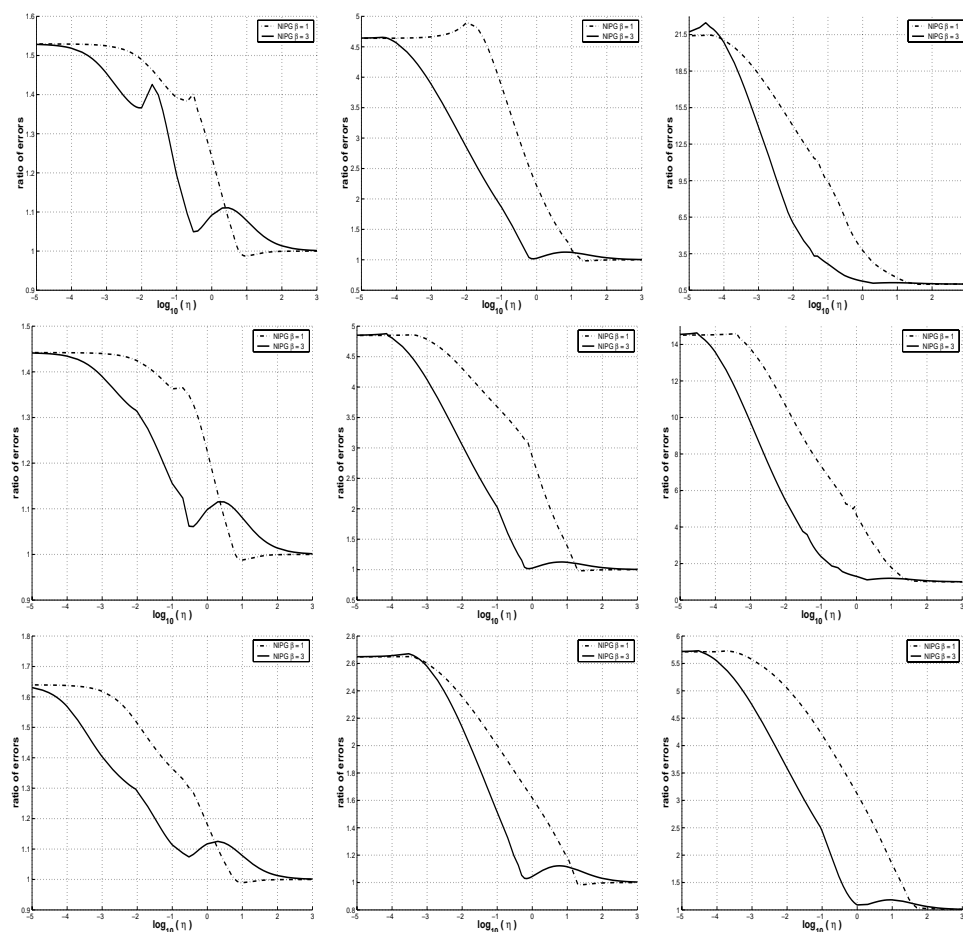
FIG. 8. *Ratio of the errors for the gradient $q_h$ NIPG : LDG. Row i corresponds to mesh i and column j to approximation of degree j.*

TABLE 4.2
*Asymptotic behavior of the spectral condition number $\kappa(h)$ as a function of the mesh size when using uniform approximations of degree p.*

| Method | Penalization | $\kappa(h)$ |
|---|---|---|
| Babuška–Zlámal | $O\left(h^{-(2p+1)}\right)$ | $O\left(h^{-(2p+2)}\right)$ |
| IP | $O\left(h^{-1}\right)$ | $O\left(h^{-2}\right)$ |
| LDG | $O\left(h^{-1}\right)$ | $O\left(h^{-2}\right)$ |
| Baumann–Oden | No penalization | $O\left(h^{-2}\right)$ |
| NIPG1 | $O\left(h^{-1}\right)$ | $O\left(h^{-2}\right)$ |
| NIPG3 | $O\left(h^{-3}\right)$ | $O\left(h^{-4}\right)$ |

FIG. 9. *Ratio of the errors for the gradient* $\nabla u_h$ *NIPG* : *LDG. Row i corresponds to mesh i and column j to approximation of degree j.*

number of the NIPG3 method is significantly larger. However, we must keep in mind the loss of accuracy of the NIPG1 method on polynomials of even degree for the $L_2$ norm.

**4.3. Effect of the stabilization parameter $\eta$.** We carry out a numerical study of the effect of the stabilization parameter $\eta$ on the quality of the approximation. Testing the accuracy of a method is a difficult if not impossible task. Ideally, for a particular mesh we would have to find the set of parameters which give the minimal error. Here, we compare the ratio of the errors between the nonsymmetric and LDG methods as a function of the stabilization parameter $\eta$ for linear, quadratic, and cubic approximations. The errors were computed using the $L_2$ norm. In Figures 7, 8, and 9, we show the ratio of the error in the potential, the gradient $q_h$, and piecewise gradient $\nabla u_h$, respectively. We have used three unstructured meshes having approximately the same number of elements: 248, 264, and 304 triangles. We observe that the behavior is independent from the distribution of the mesh points.

In Figures 10, 11, and 12, we show these ratios for meshes with different sizes: 304, 1024, and 2461 triangles. It is clear that, when $\eta$ is small, the LDG method is more
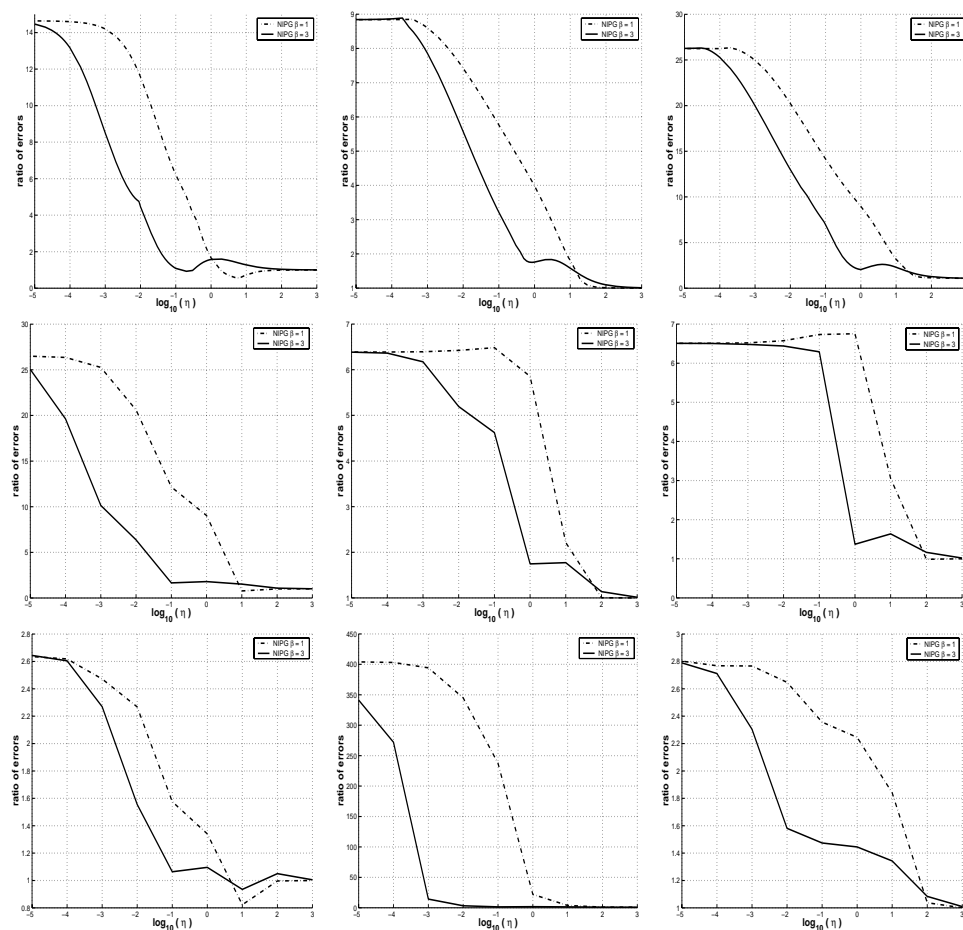
Fig. 10. *Ratio of the errors for the potential NIPG : LDG. Row i corresponds to mesh i and column j to approximation of degree j.*

accurate than both variants of the NIPG method and, consequently, the Baumann–Oden method. However, for large values of $\eta$ all the methods, including the IP (not shown in the figures), have the same accuracy.

We also compare $\boldsymbol{q_h}$ and $\nabla u_h$. We consider the ratio of the errors $R_{\nabla u}$ defined as

$$R_{\nabla u} = \frac{\|\nabla u - \nabla u_h\|_{\Omega}}{\|\nabla u - \boldsymbol{q_h}\|_{\Omega}}.$$

In Figure 13, we compare this ratio. For the IP method we considered only large values of $\eta$; for the LDG and NIPG methods we considered small values as well. In general, we observe two regimes. For $\eta \ll 1$, $\boldsymbol{q_h}$ is more accurate for the LDG method than for the NIPG methods. For $\eta \gg 1$ both approximations have the same accuracy. So in this case, the approximation obtained from the gradient operator is more efficient since it is completely local and offers the same accuracy as the auxiliary variable $\boldsymbol{q_h}$.

**5. Conclusions.** In this paper, we present the first numerical comparison of DG methods for a model elliptic problem. We give a theoretical analysis of the behavior
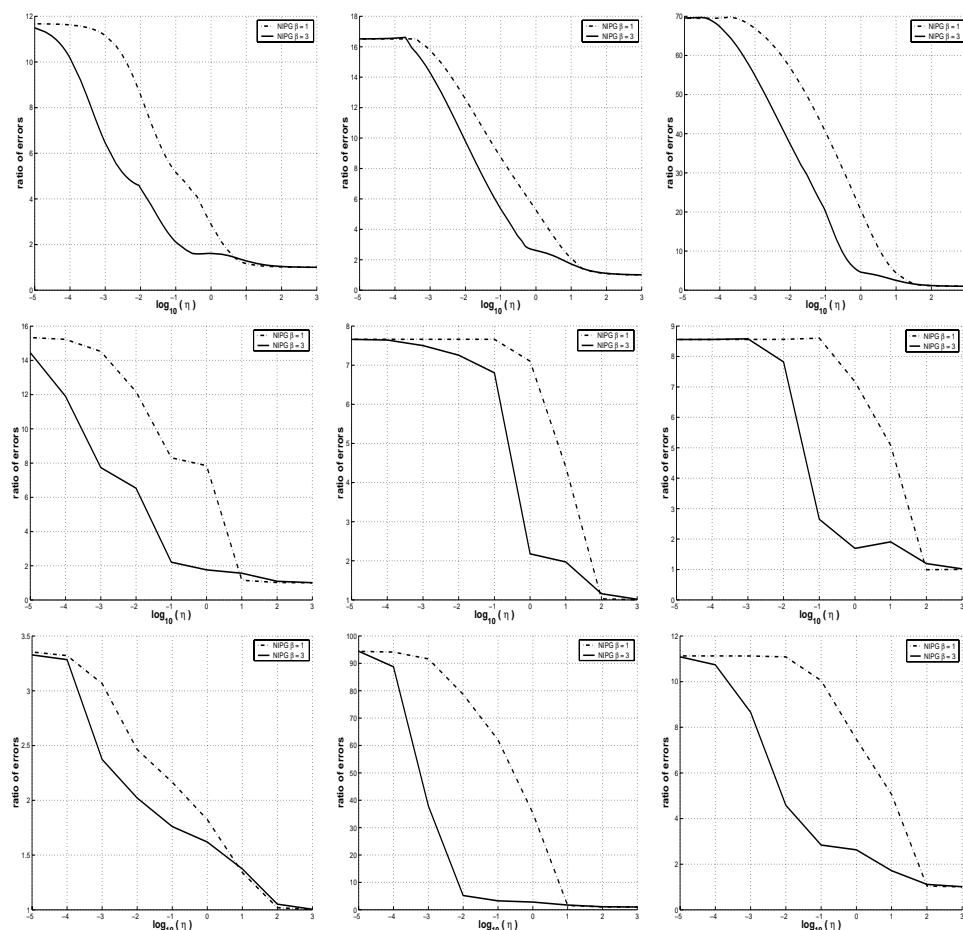
FIG. 11. *Ratio of the errors for the gradient $q_h$ NIPG : LDG. Row i corresponds to mesh i and column j to approximation of degree j.*

of the spectral condition number for methods with symmetric bilinear forms in terms of the parameters of the method, which has been shown to be sharp.

From our numerical experiments we can extract the following conclusions:

- The conditioning for the IP and LDG methods is asymptotically of the same order as for the standard continuous case.
- The nonsymmetric methods can achieve optimal rates of convergence for the potential only by using large penalty terms, hence increasing conditioning of the method from $h^{-2}$ to $h^{-4}$ when $\beta$ increases from 1 to 3. This can severely degrade the performance of the iterative method used for solving the linear system.
- In terms of efficiency in solving the linear system, methods with symmetric discretizations, i.e., LDG, IP, perform better than those with nonsymmetric discretizations, since more efficient iterative solvers such as conjugate gradient can be used.
- The LDG method requires at most 2.5 times the storage of the methods with compact stencil such as IP, Baumann–Oden, and NIPG. This has a negative
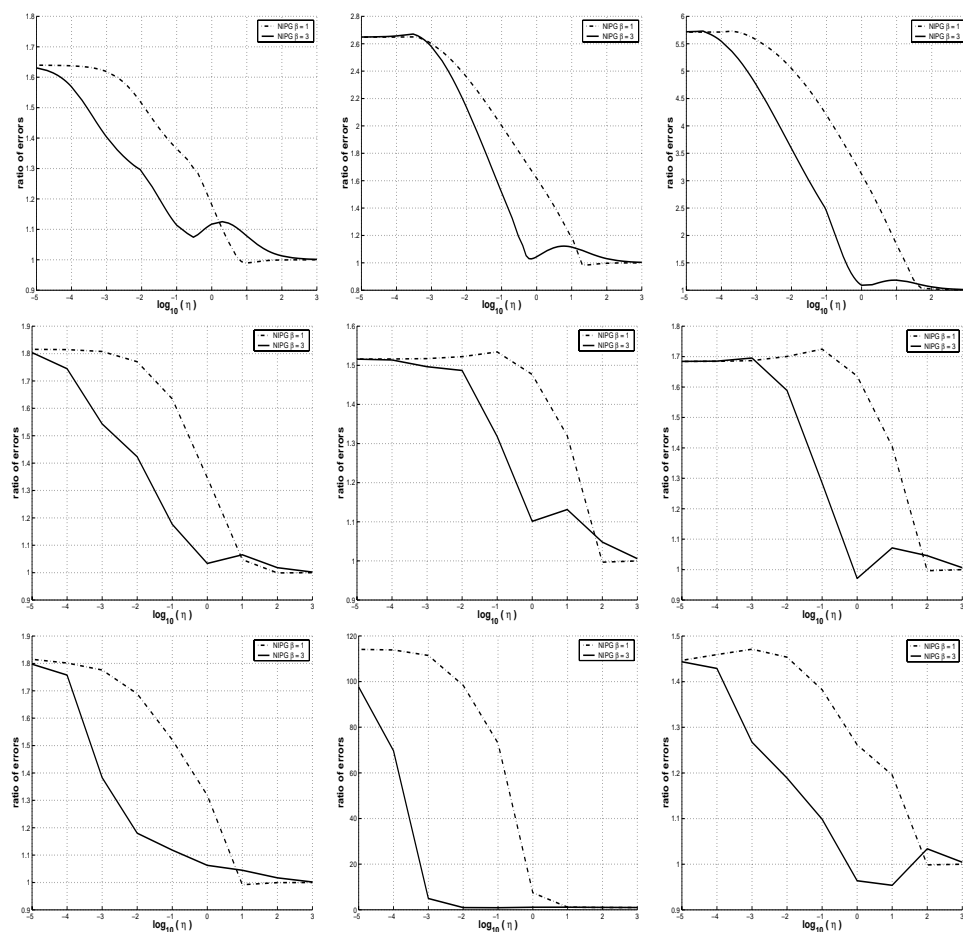
FIG. 12. *Ratio of the errors for the gradient $\nabla u_h$ NIPG : LDG. Row $i$ corresponds to mesh $i$ and column $j$ to approximation of degree $j$.*

 

       impact on the operation counts and parallel implementations.

- For larger values of $\eta$, the asymptotic behavior of the spectral condition number and accuracy of both symmetric methods, IP and LDG, are identical. In this range of penalization parameter, since the LDG method is a more expensive method, IP is more efficient.

- The LDG method is more stable than the IP method since its stabilization parameter does not depend on either the mesh or on the approximation polynomial degree.

- Although the LDG method becomes less stable as its stabilization parameter $\eta$ approaches 0, the method is more accurate than the nonsymmetric methods, such as the Baumann–Oden method.

- For large values of $\eta$, all the methods have similar accuracy in the potential and gradients $q_h$ and $\nabla u_h$. However, the gradient $q_h$ of the IP and LDG methods is more accurate than the piecewise gradient $\nabla u_h$ for small values of $\eta$. For nonsymmetric methods the opposite is true.
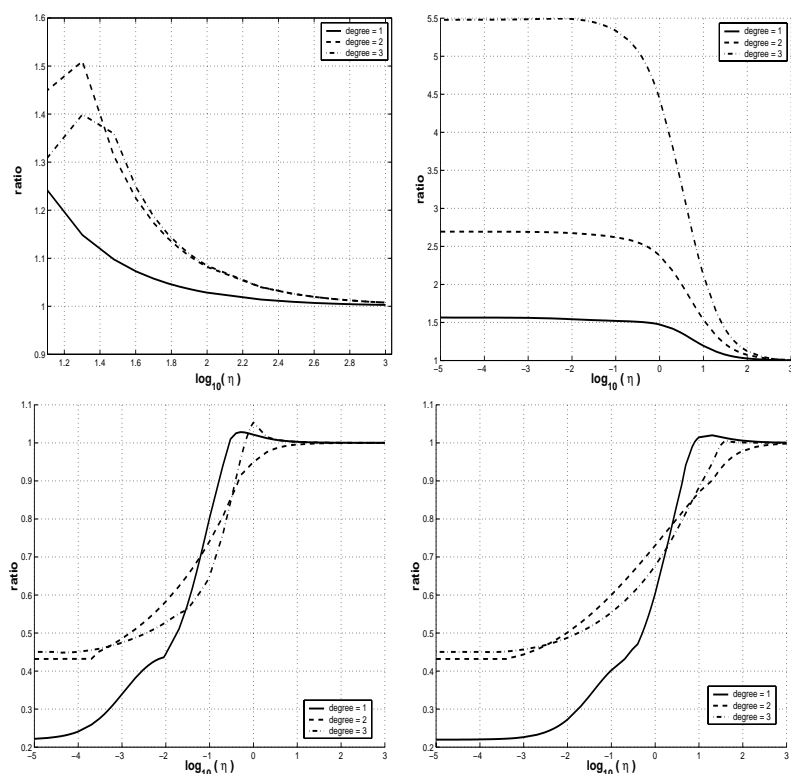
FIG. 13. *Ratio $R_{\nabla u}$ for linear, quadratic and cubic approximations. Starting from top left in clockwise order, IP, LDG, NIPG1, NIPG3.*

**Acknowledgments.** The author would like to thank Professor Bernardo Cockburn and Dr. Ilaria Perugia for insightful discussions and thanks the anonymous reviewers for their valuable suggestions.

## REFERENCES

[1]  S. AGMON, *Lectures on Elliptic Boundary Value Problems*, Van Nostrand, Princeton, NJ, 1965.
[2]  D. N. ARNOLD, *An interior penalty finite element method with discontinuous elements*, SIAM J. Numer. Anal., 19 (1982), pp. 742–760.
[3]  D. ARNOLD, F. BREZZI, B. COCKBURN, AND D. MARINI, *DG methods for elliptic problems*, in Discontinuous Galerkin Methods: Theory, Computation and Applications, B. Cockburn, G. Karniadakis, and C. Shu, eds., Lecture Notes in Comput. Sci. Engrg. 11, Springer-Verlag, Berlin, 2000, pp. 89–101.
[4]  I. BABUŠKA, C. BAUMANN, AND J. ODEN, *A discontinuous hp finite element method for diffusion problems: 1-D analysis*, Comput. Math. Appl., 37 (1999), pp. 103–122.
[5]  I. BABUŠKA AND M. ZLÁMAL, *Nonconforming elements in the finite element method with penalty*, SIAM J. Numer. Anal., 10 (1973), pp. 863–875.
[6]  G. BAKER, *Finite element methods for elliptic equations using nonconforming elements*, Math. Comp., 31 (1977), pp. 45–59.
[7]  T. BARTH AND H. DECONINCK, EDS., *High-Order Methods for Computational Physics*, Lecture Notes in Comput. Sci. Engrg. 9, Springer-Verlag, Berlin, 1999.
[8]  F. BASSI AND S. REBAY, *A high-order accurate discontinuous finite element method for the numerical solution of the compressible Navier-Stokes equations*, J. Comput. Phys., 131 (1997), pp. 267–279.

[9]   C. Baumann and J. Oden, *A discontinuous hp finite element method for convection-diffusion problems*, Comput. Methods Appl. Mech. Engrg., 175 (1999), pp. 311–341.

[10]  P. Castillo, *Local Discontinuous Galerkin Methods for Convection-Diffusion and Elliptic Problems*, Ph.D. thesis, University of Minnesota, Minneapolis, MN, 2001.

[11]  P. Castillo, B. Cockburn, I. Perugia, and D. Schötzau, *An a priori error analysis of the local discontinuous Galerkin method for elliptic problems*, SIAM J. Numer. Anal., 38 (2000), pp. 1676–1706.

[12]  B. Cockburn and C. Dawson, *Some extensions of the local discontinuous Galerkin method for convection-diffusion equations in multidimensions*, in Proceedings of the Conference on the Mathematics of Finite Elements and Applications: MAFELAP X, J. Whiteman, ed., Elsevier, New York, 2000, pp. 225–238.

[13]  B. Cockburn, G. Kanschat, I. Perugia, and D. Schötzau, *Superconvergence of the local discontinuous Galerkin method for elliptic problems on Cartesian grids*, SIAM J. Numer. Anal., 39 (2001), pp. 264–285.

[14]  B. Cockburn, G. Karniadakis, and C. Shu, eds., *Discontinuous Galerkin Methods: Theory, Computation and Applications*, Lecture Notes in Comput. Sci. Engrg. 11, Springer-Verlag, Berlin, 2000.

[15]  B. Cockburn and C.-W. Shu, *The local discontinuous Galerkin method for convection-diffusion systems*, SIAM J. Numer. Anal., 35 (1998), pp. 2440–2463.

[16]  J. Douglas, Jr. and T. Dupont, *Interior Penalty Procedures for Elliptic and Parabolic Galerkin Methods*, Lecture Notes in Phys. 58, Springer-Verlag, Berlin, 1976.

[17]  M. Dubiner, *Spectral methods on triangles and other domains*, J. Sci. Comput., 6 (1991), pp. 345–390.

[18]  L. C. Evans, *Partial Differential Equations*, Graduate Studies in Mathematics 19, AMS, Providence, RI, 1998.

[19]  J. Oden, I. Babuška, and C. Baumann, *A discontinuous hp finite element method for diffusion problems*, J. Comput. Phys., 146 (1998), pp. 491–519.

[20]  A. Quarteroni and A. Valli, *Numerical Approximation of Partial Differential Equations*, Springer Ser. Comput. Math. 23, Springer-Verlag, Berlin, 1994.

[21]  B. Rivière, M. F. Wheeler, and V. Girault, *Improved energy estimates for interior penalty, constrained and discontinuous Galerkin methods for elliptic problems. Part* I, Comput. Geosci., 3 (1999), pp. 337–360.

[22]  B. Rivière, M. F. Wheeler, and V. Girault, *A priori error estimates for finite element methods based on discontinuous approximation spaces for elliptic problems*, SIAM J. Numer. Anal., 39 (2001), pp. 902–931.

[23]  E. Süli, C. Schwab, and P. Houston, *hp-DGFEM for partial differential equations with nonnegative characteristic form*, in Discontinuous Galerkin Methods: Theory, Computation and Applications, B. Cockburn, G. Karniadakis, and C. Shu, eds., Lecture Notes in Comput. Sci. Engrg. 11, Springer-Verlag, Berlin, 2000, pp. 221–230.

[24]  M. F. Wheeler, *An elliptic collocation-finite element method with interior penalties*, SIAM J. Numer. Anal., 15 (1978), pp. 152–161.