

Chapter 2

DGM for Elliptic Problems

This chapter concerns in basic aspects of the discontinuous Galerkin method (DGM), which will be treated in an example of a simple problem for the Poisson equation with mixed Dirichlet–Neumann boundary conditions. We introduce the discretization of this problem with the aid of several variants of the DGM. Further, we prove the existence of the approximate solution and derive error estimates. Finally, several numerical examples are presented.

The book contains a detailed analysis of qualitative properties of DG techniques. It is based on a number of estimates with various constants. We denote by $C_A, C_B, C_C, \dots, C_a, C_b, C_c, \dots$ positive constants arising in the formulation of results that can be simply named (e.g., $_A$ corresponds to approximation properties, $_B$ —boundedness, $_C$ —coercivity, etc.) Otherwise, we use symbols C, C_1, C_2, \dots . These constants are always independent of the parameters of the discretization (i.e., the space mesh-size h , time step τ in the case of nonstationary problems, and also the degree p of polynomial approximation in the case of the hp -methods), but they may depend on the data in problems. They are often “autonomous” in individual chapters or sections. Some constants are sometimes defined in a complicated way on the basis of a number of constants appearing in previous considerations. For an example, see Remark 4.13.

2.1 Model Problem

Let Ω be a bounded domain in \mathbb{R}^d , $d = 2, 3$, with Lipschitz boundary $\partial\Omega$. We denote by $\partial\Omega_D$ and $\partial\Omega_N$ parts of the boundary $\partial\Omega$ such that $\partial\Omega = \partial\Omega_D \cup \partial\Omega_N$, $\partial\Omega_D \cap \partial\Omega_N = \emptyset$ and $\partial\Omega_D \neq \emptyset$.

We consider the following model problem for the Poisson equation: Find a function $u : \Omega \rightarrow \mathbb{R}$ such that

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

$$u = u_D \quad \text{on } \partial\Omega_D, \quad (2.1b)$$

$$\mathbf{n} \cdot \nabla u = g_N \quad \text{on } \partial\Omega_N, \quad (2.1c)$$

where f , u_D and g_N are given functions. Let us note that $\mathbf{n} \cdot \nabla u = \frac{\partial u}{\partial \mathbf{n}}$ is the derivative of the function u in the direction \mathbf{n} , which is the outer unit normal to $\partial\Omega$. A function $u \in C^2(\overline{\Omega})$ satisfying (2.1) pointwise is called a *classical solution*. It is suitable to introduce a weak formulation of the above problem. Let us define the space

$$V = \{v \in H^1(\Omega); v|_{\partial\Omega_D} = 0\}.$$

Assuming that u is a classical solution, we multiply (2.1a) by any function $v \in V$, integrate over Ω and use Green's theorem. Taking into account the boundary condition (2.1c), we obtain the identity

$$\int_{\Omega} \nabla u \cdot \nabla v \, dx = \int_{\Omega} f v \, dx + \int_{\partial\Omega_N} g_N v \, dS \quad \forall v \in V. \quad (2.2)$$

We can introduce the following definition.

Definition 2.1 Let us assume the existence of $u^* \in H^1(\Omega)$ such that $u^*|_{\partial\Omega_D} = u_D$ and let $f \in L^2(\Omega)$, $g_N \in L^2(\partial\Omega_N)$. Now we say that a function u is a *weak solution* of problem (2.1), if

- (a) $u - u^* \in V$,
- (b) u satisfies identity (2.2).

Using the Lax–Milgram Lemma 1.6, we can prove that there exists a unique weak solution of (2.1), see, e.g., [233, Sect. 6.1.2]. In the following, we deal with numerical solution of problem (2.1) with the aid of discontinuous piecewise polynomial approximations.

2.2 Abstract Numerical Method and Its Theoretical Analysis

In order to better understand theoretical foundations of the DGM, we describe a possible general approach to deriving error estimates. (Readers familiar with concepts of a priori error estimates in the finite element method can skip this section.)

Let $u \in V$ be a weak solution of a given problem. Let V_h denote a *finite-dimensional space*, where an *approximate solution* u_h is sought. The subscript $h > 0$ (usually chosen as $h \in (0, \bar{h})$ with $\bar{h} > 0$) denotes the parameter of the discretization. Further, we introduce an infinitely dimensional function space W_h such that $V \subset W_h$ and $V_h \subset W_h$. (If $V_h \subset V$, then we usually put $W_h := V$ and thus, W_h is independent of h .) Finally, let $\|\cdot\|_{W_h}$ be a suitable norm in W_h . As we see later,

the spaces V_h and W_h will be constructed over a suitable mesh in the computational domain, and hence the norm $\|\cdot\|_{W_h}$ may be mesh-dependent.

An *abstract numerical method* reads: Find $u_h \in V_h$ such that

$$A_h(u_h, v_h) = F(v_h) \quad \forall v_h \in V_h, \quad (2.3)$$

where $A_h : W_h \times W_h \rightarrow \mathbb{R}$ is a bilinear form and $F : W_h \rightarrow \mathbb{R}$ is a linear functional.

In the numerical analysis, we want to reach the following goals:

- the approximate solution u_h of (2.3) *exists* and is *unique*,
- the approximate solution u_h *converges* to the exact solution u in the $\|\cdot\|_{W_h}$ -norm as $h \rightarrow 0$, i.e.,

$$\lim_{h \rightarrow 0} \|u - u_h\|_{W_h} = 0, \quad (2.4)$$

- *a priori error estimate*, i.e., we seek $\alpha > 0$ independent of h such that

$$\|u - u_h\|_{W_h} \leq Ch^\alpha, \quad h \in (0, \bar{h}), \quad (2.5)$$

where $C > 0$ is a constant, independent of h (but may depend on u), and α is the *order of convergence*.

Obviously, an *a priori* error estimate implies the convergence.

The existence and uniqueness of the approximate solution is a consequence of the *coercivity* of A_h , i.e., there exists $C_c > 0$ such that

$$A_h(v_h, v_h) \geq C_c \|v_h\|_{W_h}^2 \quad \forall v_h \in V_h. \quad (2.6)$$

Then Corollary 1.7 implies the existence and uniqueness of the approximate solution u_h .

In order to derive *a priori* error estimates, we prove the *consistency* of the method,

$$A_h(u, v_h) = F(v_h) \quad \forall v_h \in V_h \quad (2.7)$$

which, together with (2.3), immediately gives the *Galerkin orthogonality* of the error $e_h = u_h - u$ to the space V_h :

$$A_h(e_h, v_h) = 0 \quad \forall v_h \in V_h. \quad (2.8)$$

Further, we introduce an *interpolation operator* (usually defined as a suitable *projection*) $\Pi_h : V \rightarrow V_h$ and prove its *approximation property*, namely existence of a constant $\alpha > 0$ such that

$$\|v - \Pi_h v\|_{W_h} \leq \tilde{C}(v)h^\alpha \quad \forall v \in V, \quad h \in (0, \bar{h}), \quad (2.9)$$

where $\tilde{C}(v) > 0$ is a constant independent of h but dependent on v . A further step is the derivation of the inequality

$$A_h(u - \Pi_h u, v_h) \leq R(u - \Pi_h u) \|v_h\|_{W_h} \quad \forall v_h \in V_h, \quad (2.10)$$

where R depends on suitable norms of the interpolation error $u - \Pi_h u$.

Finally, the *error estimate* is derived in the following way: for each $h \in (0, \bar{h})$ we decompose the error e_h by

$$e_h = u_h - u = \xi + \eta, \quad (2.11)$$

where $\xi := u_h - \Pi_h u \in V_h$ and $\eta := \Pi_h u - u \in W_h$. Putting $v_h := \xi$ in (2.8), we get

$$A_h(e_h, \xi) = A_h(\xi, \xi) + A_h(\eta, \xi) = 0. \quad (2.12)$$

It follows from the coercivity (2.6) and estimate (2.10) that

$$C_c \|\xi\|_{W_h}^2 \leq A_h(\xi, \xi) = -A_h(\eta, \xi) \leq R(\eta) \|\xi\|_{W_h}, \quad (2.13)$$

which immediately implies the inequality

$$\|\xi\|_{W_h} \leq \frac{R(\eta)}{C_c}. \quad (2.14)$$

Now, the triangle inequality, relations (2.11) and (2.14) give the error estimate in the form

$$\|e_h\|_{W_h} \leq \|\xi\|_{W_h} + \|\eta\|_{W_h} \leq \frac{R(\eta)}{C_c} + \|\eta\|_{W_h}. \quad (2.15)$$

This is often called the *abstract error estimate*, which represents an error bound in terms of the interpolation error η .

The last aim is to use the approximation property (2.9) of the operator Π_h and to estimate the expression $R(\eta)$ in terms of the mesh-size h in the form

$$R(\eta) \leq \tilde{C}_1(u) h^\alpha, \quad (2.16)$$

which together with (2.15) immediately imply the *error estimate*

$$\|e_h\|_{W_h} \leq \left(C_c^{-1} \tilde{C}_1(u) + \tilde{C}(u) \right) h^\alpha, \quad (2.17)$$

valid for all $h \in (0, \bar{h})$. We say that the numerical scheme has the *order of convergence* in the norm $\|\cdot\|_{W_h}$ equal to α .

This concept of numerical analysis is applied in this chapter. (Among other, we specify there the spaces W_h and V_h .) For time dependent problems, treated in Chaps. 4–6, the analysis is more complicated and the previous technique has to be modified. However, in some parts of the book, error estimates are derived in a different way.

Remark 2.2 As was mentioned above, we are interested here in deriving of *a priori error estimates* (simply called *error estimates*). We do not deal with *a posteriori error estimates*, when the error is bounded in a suitable norm in terms of the approximate solution and data of the problem. The subject of *a posteriori error estimates* plays an important role in practical computations, but is out of the scope of this book. For some results in this direction for the DGM we can refer, e.g., to the papers [5, 91, 118, 166, 185, 190] and the references cited therein.

2.3 Spaces of Discontinuous Functions

The subject of this section is the construction of DG space partitions of the bounded computational domain Ω and the specification of their properties which are used in the theoretical analysis. Further, function spaces over these meshes are defined.

2.3.1 Partition of the Domain

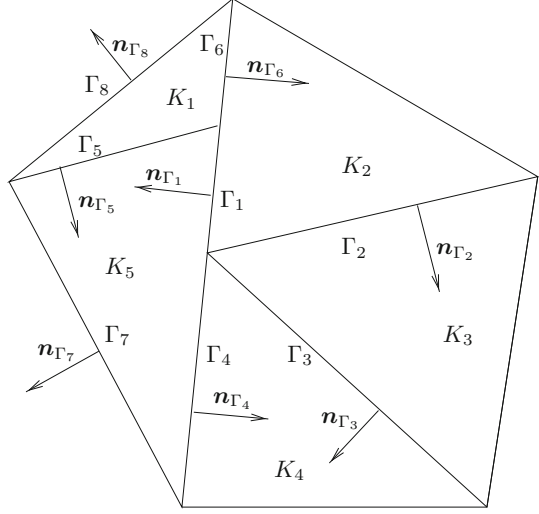
Let \mathcal{T}_h ($h > 0$ is a parameter) be a partition of the closure $\overline{\Omega}$ of the domain Ω into a finite number of closed d -dimensional simplexes K with mutually disjoint interiors such that

$$\overline{\Omega} = \bigcup_{K \in \mathcal{T}_h} K. \quad (2.18)$$

This assumption means that the domain Ω is polygonal (if $d = 2$) or polyhedral (if $d = 3$). The case of a 2D nonpolygonal domain is considered, e.g., in [256], where curved elements are used. See also Chap. 8, where curved elements are treated from the implementation point of view. We call \mathcal{T}_h a *triangulation* of Ω and do not require the standard conforming properties from the finite element method, introduced e.g., in [37, 52, 115, 254] or [287]. In two-dimensional problems ($d = 2$) we choose $K \in \mathcal{T}_h$ as triangles and in three-dimensional problems ($d = 3$) the elements $K \in \mathcal{T}_h$ are tetrahedra. As we see, we admit that in the finite element mesh the so-called *hanging nodes* (and in 3D also *hanging edges*) appear; see Fig. 2.1.

In general, the discontinuous Galerkin method can handle with more general elements as quadrilaterals and convex or even nonconvex star-shaped polygons in 2D and hexahedra, pyramids and convex or nonconvex star-shaped polyhedra in 3D.

Fig. 2.1 Example of elements K_l , $l = 1, \dots, 5$, and faces Γ_l , $l = 1, \dots, 8$, with the corresponding normals \mathbf{n}_{Γ_l} . The triangle K_5 has a hanging node. Its boundary is formed by four edges: $\partial K_5 = \Gamma_1 \cup \Gamma_4 \cup \Gamma_7 \cup \Gamma_5$



As an example, we can consider the so-called dual finite volumes constructed over triangular ($d = 2$) or tetrahedral ($d = 3$) meshes (cf., e.g., [126]). A use of such elements will be discussed in Sect. 7.2.

In our further considerations we use the following notation. By ∂K we denote the boundary of an element $K \in \mathcal{T}_h$ and set $h_K = \text{diam}(K) = \text{diameter of } K$, $h = \max_{K \in \mathcal{T}_h} h_K$. By ρ_K we denote the radius of the largest d -dimensional ball inscribed into K and by $|K|$ we denote the d -dimensional Lebesgue measure of K .

Let $K, K' \in \mathcal{T}_h$. We say that K and K' are *neighbouring elements* (or simply *neighbours*) if the set $\partial K \cap \partial K'$ has positive $(d - 1)$ -dimensional measure. We say that $\Gamma \subset K$ is a *face* of K , if it is a maximal connected open subset of either $\partial K \cap \partial K'$, where K' is a neighbour of K , or $\partial K \cap \partial \Omega_D$ or $\partial K \cap \partial \Omega_N$. The symbol $|\Gamma|$ will denote the $(d - 1)$ -dimensional Lebesgue measure of Γ . Hence, if $d = 2$, then $|\Gamma|$ is the length of Γ and for $d = 3$, $|\Gamma|$ denotes the area of Γ . By \mathcal{F}_h we denote the system of all faces of all elements $K \in \mathcal{T}_h$. Further, we define the set of all boundary faces by

$$\mathcal{F}_h^B = \{\Gamma \in \mathcal{F}_h; \Gamma \subset \partial \Omega\},$$

the set of all “Dirichlet” boundary faces by

$$\mathcal{F}_h^D = \{\Gamma \in \mathcal{F}_h; \Gamma \subset \partial \Omega_D\},$$

the set of all “Neumann” boundary faces by

$$\mathcal{F}_h^N = \{\Gamma \in \mathcal{F}_h, \Gamma \subset \partial \Omega_N\}$$

and the set of all inner faces

$$\mathcal{F}_h^I = \mathcal{F}_h \setminus \mathcal{F}_h^B.$$

Obviously, $\mathcal{F}_h = \mathcal{F}_h^I \cup \mathcal{F}_h^D \cup \mathcal{F}_h^N$ and $\mathcal{F}_h^B = \mathcal{F}_h^D \cup \mathcal{F}_h^N$. For a shorter notation we put

$$\mathcal{F}_h^{ID} = \mathcal{F}_h^I \cup \mathcal{F}_h^D.$$

For each $\Gamma \in \mathcal{F}_h$ we define a unit normal vector \mathbf{n}_Γ . We assume that for $\Gamma \in \mathcal{F}_h^B$ the normal \mathbf{n}_Γ has the same orientation as the outer normal to $\partial\Omega$. For each face $\Gamma \in \mathcal{F}_h^I$ the orientation of \mathbf{n}_Γ is arbitrary but fixed. See Fig. 2.1.

For each $\Gamma \in \mathcal{F}_h^I$ there exist two neighbouring elements $K_\Gamma^{(L)}, K_\Gamma^{(R)} \in \mathcal{T}_h$ such that $\Gamma \subset \partial K_\Gamma^{(L)} \cap \partial K_\Gamma^{(R)}$. (This means that the elements $K_\Gamma^{(L)}, K_\Gamma^{(R)}$ are adjacent to Γ and they share this face.) We use the convention that \mathbf{n}_Γ is the outer normal to $\partial K_\Gamma^{(L)}$ and the inner normal to $\partial K_\Gamma^{(R)}$; see Fig. 2.2.

Moreover, if $\Gamma \in \mathcal{F}_h^B$, then there exists an element $K_\Gamma^{(L)} \in \mathcal{T}_h$ such that $\Gamma \subset K_\Gamma^{(L)} \cap \partial\Omega$.

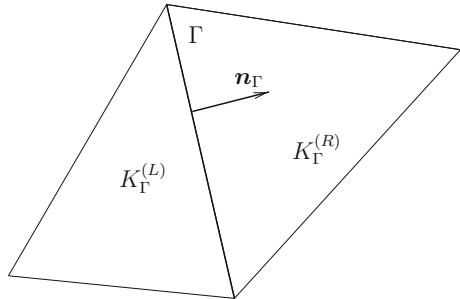
2.3.2 Assumptions on Meshes

Let us consider a system $\{\mathcal{T}_h\}_{h \in (0, \bar{h})}$, $\bar{h} > 0$, of triangulations of the domain Ω ($\mathcal{T}_h = \{K\}_{K \in \mathcal{T}_h}$). In our further considerations we meet various assumptions on triangulations. The first is usual in the theory of the finite element method:

- The system $\{\mathcal{T}_h\}_{h \in (0, \bar{h})}$ of triangulations is *shape-regular*: there exists a positive constant C_R such that

$$\frac{h_K}{\rho_K} \leq C_R \quad \forall K \in \mathcal{T}_h \quad \forall h \in (0, \bar{h}). \quad (2.19)$$

Fig. 2.2 Interior face Γ , elements $K_\Gamma^{(L)}$ and $K_\Gamma^{(R)}$ and the orientation of \mathbf{n}_Γ



Moreover, for each face $\Gamma \in \mathcal{F}_h$, $h \in (0, \bar{h})$, we need to introduce a quantity $h_\Gamma > 0$, which represents a “one-dimensional” size of the face Γ . We require that

- the quantity h_Γ satisfies the *equivalence condition* with h_K , i.e., there exist constants $C_T, C_G > 0$ independent of h, K and Γ such that

$$C_T h_K \leq h_\Gamma \leq C_G h_K, \quad \forall K \in \mathcal{T}_h, \forall \Gamma \in \mathcal{F}_h, \Gamma \subset \partial K, \forall h \in (0, \bar{h}). \quad (2.20)$$

The equivalence condition can be fulfilled by additional assumptions on the system of triangulations $\{\mathcal{T}_h\}_{h \in (0, \bar{h})}$ and by a suitable choice of the quantity h_Γ , $\Gamma \in \mathcal{F}_h$, $h \in (0, \bar{h})$. We introduce some assumptions on triangulations and several choices of the quantity h_Γ . Then we discuss how the equivalence condition (2.20) is satisfied.

In literature we can find the following assumptions on the system of triangulations:

(MA1) The system $\{\mathcal{T}_h\}_{h \in (0, \bar{h})}$ is *locally quasi-uniform*: there exists a constant $C_Q > 0$ such that

$$h_K \leq C_Q h_{K'}, \quad \forall K, K' \in \mathcal{T}_h, K, K' \text{ are neighbours}, \forall h \in (0, \bar{h}). \quad (2.21)$$

(MA2) The faces $\Gamma \subset \partial K$ do not degenerate with respect to the diameter of K if $h \rightarrow 0$: there exists a constant $C_d > 0$ such that

$$h_K \leq C_d \text{diam}(\Gamma) \quad \forall K \in \mathcal{T}_h \quad \forall \Gamma \in \mathcal{F}_h, \Gamma \subset \partial K, \quad \forall h \in (0, \bar{h}). \quad (2.22)$$

(MA3) The system $\{\mathcal{T}_h\}_{h \in (0, \bar{h})}$ is *quasi-uniform*: there exists a constant $C_U > 0$ such that

$$h \leq C_U h_K \quad \forall K \in \mathcal{T}_h \quad \forall h \in (0, \bar{h}). \quad (2.23)$$

(MA4) The triangulations \mathcal{T}_h , $h \in (0, \bar{h})$, are *conforming*. This means that for two elements $K, K' \in \mathcal{T}_h$, $K \neq K'$, either $K \cap K' = \emptyset$ or $K \cap K'$ is a common vertex or $K \cap K'$ is a common face (or for $d = 3$, when $K \cap K'$ is a common edge) of K and K' .

If condition (MA4) is not satisfied, then the triangulations \mathcal{T}_h are called *nonconforming*.

Remark 2.3 There are some relations among the mesh assumptions (MA1)–(MA4) mentioned above. Obviously, (MA3) \Rightarrow (MA1). Moreover, if the system of triangulation is shape-regular (i.e., (2.19) is fulfilled) then (MA4) \Rightarrow (MA1) & (MA2).

Exercises 2.4 Prove the implications in Remark 2.3.

Concerning the choice of the quantity h_Γ , $\Gamma \in \mathcal{F}_h$, $h \in (0, \bar{h})$, in literature we can find the following basic possibilities:

$$(i) \quad h_\Gamma = \text{diam}(\Gamma), \quad \Gamma \in \mathcal{T}_h, \quad (2.24)$$

$$(ii) \quad h_\Gamma = \begin{cases} \frac{1}{2} \left(h_{K_\Gamma^{(L)}} + h_{K_\Gamma^{(R)}} \right) & \text{for } \Gamma \in \mathcal{F}_h^I \\ h_{K_\Gamma^{(L)}} & \text{for } \Gamma \in \mathcal{F}_h^B, \end{cases} \quad (2.25)$$

$$(iii) \quad h_\Gamma = \begin{cases} \max \left(h_{K_\Gamma^{(L)}}, h_{K_\Gamma^{(R)}} \right) & \text{for } \Gamma \in \mathcal{F}_h^I \\ h_{K_\Gamma^{(L)}} & \text{for } \Gamma \in \mathcal{F}_h^B, \end{cases} \quad (2.26)$$

$$(iv) \quad h_\Gamma = \begin{cases} \min \left(h_{K_\Gamma^{(L)}}, h_{K_\Gamma^{(R)}} \right) & \text{for } \Gamma \in \mathcal{F}_h^I \\ h_{K_\Gamma^{(L)}} & \text{for } \Gamma \in \mathcal{F}_h^B, \end{cases} \quad (2.27)$$

where $K_\Gamma^{(L)}, K_\Gamma^{(R)} \in \mathcal{T}_h$ are the elements adjacent to $\Gamma \in \mathcal{F}_h^I$, see Fig. 2.2, and $K_\Gamma^{(L)} \in \mathcal{T}_h$ is the element adjacent to $\Gamma \in \mathcal{F}_h^B$.

The following lemma characterizes assumptions on computational grids and the choice of h_Γ , which guarantee the equivalence condition (2.20).

Lemma 2.5 *Let $\{\mathcal{T}_h\}_{h \in (0, \bar{h})}$ be a system of triangulations of the domain Ω satisfying the shape-regularity assumption (2.19). Then the equivalence condition (2.20) is satisfied in the following cases:*

- (i) *The triangulations \mathcal{T}_h , $h \in (0, \bar{h})$, are conforming (i.e., assumption (MA4) is satisfied) and h_Γ are defined by (2.24) or (2.25) or (2.26) or (2.27).*
- (ii) *The triangulations \mathcal{T}_h , $h \in (0, \bar{h})$, are, in general, nonconforming; assumption (MA2) (i.e., (2.22)) is satisfied and h_Γ are defined by (2.24).*
- (iii) *The triangulations \mathcal{T}_h , $h \in (0, \bar{h})$, are, in general, nonconforming; assumption (MA1) is satisfied (i.e., the system $\{\mathcal{T}_h\}_{h \in (0, \bar{h})}$ is locally quasi-uniform) and h_Γ are defined by (2.25) or (2.26) or (2.27).*

Exercises 2.6 Prove the above lemma and find the constants C_T and C_G . For example, in the case (iii), when h_Γ is given by (2.25), we have

$$C_T = (1 + C_Q^{-1})/2, \quad C_G = (1 + C_Q)/2, \quad (2.28)$$

where C_Q is the constant from the local quasi-uniformity condition (2.21).

2.3.3 Broken Sobolev Spaces

The discontinuous Galerkin method is based on the use of discontinuous approximations. This is the reason that over a triangulation \mathcal{T}_h , for any $k \in \mathbb{N}$, we define the so-called *broken Sobolev space*

$$H^k(\Omega, \mathcal{T}_h) = \{v \in L^2(\Omega); v|_K \in H^k(K) \forall K \in \mathcal{T}_h\}, \quad (2.29)$$

which consists of functions, whose restrictions on $K \in \mathcal{T}_h$ belong to the Sobolev space $H^k(K)$. On the other hand, functions from $H^k(\Omega, \mathcal{T}_h)$ are, in general, discontinuous on inner faces of elements $K \in \mathcal{T}_h$. For $v \in H^k(\Omega, \mathcal{T}_h)$, we define the norm

$$\|v\|_{H^k(\Omega, \mathcal{T}_h)} = \left(\sum_{K \in \mathcal{T}_h} \|v\|_{H^k(K)}^2 \right)^{1/2} \quad (2.30)$$

and the seminorm

$$|v|_{H^k(\Omega, \mathcal{T}_h)} = \left(\sum_{K \in \mathcal{T}_h} |v|_{H^k(K)}^2 \right)^{1/2}. \quad (2.31)$$

Let $\Gamma \in \mathcal{F}_h^I$ and let $K_\Gamma^{(L)}, K_\Gamma^{(R)} \in \mathcal{T}_h$ be elements adjacent to Γ . For $v \in H^1(\Omega, \mathcal{T}_h)$ we introduce the following notation:

$$\begin{aligned} v_\Gamma^{(L)} &= \text{the trace of } v|_{K_\Gamma^{(L)}} \text{ on } \Gamma, \\ v_\Gamma^{(R)} &= \text{the trace of } v|_{K_\Gamma^{(R)}} \text{ on } \Gamma, \\ \langle v \rangle_\Gamma &= \frac{1}{2} \left(v_\Gamma^{(L)} + v_\Gamma^{(R)} \right) \quad (\text{mean value of the traces of } v \text{ on } \Gamma), \\ [v]_\Gamma &= v_\Gamma^{(L)} - v_\Gamma^{(R)} \quad (\text{jump of } v \text{ on } \Gamma). \end{aligned} \quad (2.32)$$

The value $[v]_\Gamma$ depends on the orientation of \mathbf{n}_Γ , but $[v]_\Gamma \mathbf{n}_\Gamma$ is independent of this orientation.

Moreover, let $\Gamma \in \mathcal{F}_h^B$ and $K_\Gamma^{(L)} \in \mathcal{T}_h$ be the element such that $\Gamma \subset \partial K_\Gamma^{(L)} \cap \partial \Omega$. Then for $v \in H^1(\Omega, \mathcal{T}_h)$ we introduce the following notation:

$$\begin{aligned} v_\Gamma^{(L)} &= \text{the trace of } v|_{K_\Gamma^{(L)}} \text{ on } \Gamma, \\ \langle v \rangle_\Gamma &= [v]_\Gamma = v_\Gamma^{(L)}. \end{aligned} \quad (2.33)$$

If $\Gamma \in \mathcal{F}_h^B$, then by $v_\Gamma^{(R)}$ we formally denote the exterior trace of v on Γ given either by a boundary condition or by an extrapolation from the interior of Ω .

In case that $\Gamma \in \mathcal{F}_h$ and $[\cdot]_\Gamma, \langle \cdot \rangle_\Gamma$ and \mathbf{n}_Γ appear in integrals $\int_\Gamma \dots dS$, then we usually omit the subscript Γ and simply write $[\cdot], \langle \cdot \rangle$ and \mathbf{n} , respectively.

The discontinuous Galerkin method can be characterized as a finite element technique using piecewise polynomial approximations, in general discontinuous on interfaces between neighbouring elements. Therefore, we introduce a finite-dimensional subspace of $H^k(\Omega, \mathcal{T}_h)$, where the approximate solution will be sought.

Let \mathcal{T}_h be a triangulation of Ω introduced in Sect. 2.3.1 and let $p \geq 0$ be an integer. We define the space of discontinuous piecewise polynomial functions

$$S_{hp} = \{v \in L^2(\Omega); v|_K \in P_p(K) \forall K \in \mathcal{T}_h\}, \quad (2.34)$$

where $P_p(K)$ denotes the space of all polynomials of degree $\leq p$ on K . We call the number p the *degree of polynomial approximation*. Obviously, $S_{hp} \subset H^k(\Omega, \mathcal{T}_h)$ for any $k \geq 1$ and its dimension $\dim S_{hp} < \infty$.

2.4 DGM Based on a Primal Formulation

In this section we introduce the so-called discontinuous Galerkin method (DGM) based on a *primal formulation* for the solution of problem (2.1). The approximate solution will be sought in the space $S_{hp} \subset H^1(\Omega, \mathcal{T}_h)$. In contrast to the standard (conforming) finite element method, the weak formulation (2.2) given in Sect. 2.1 is not suitable for the derivation of the DGM, because (2.2) does not make sense for $u \in H^1(\Omega, \mathcal{T}_h) \not\subset H^1(\Omega)$. Therefore, we introduce a “weak form of (2.1) in the sense of broken Sobolev spaces”.

Let us assume that u is a sufficiently regular solution of (2.1), namely, let $u \in H^2(\Omega)$. Then we speak of a *strong solution*. In deriving the DGM we proceed in the following way. We multiply (2.1a) by a function $v \in H^1(\Omega, \mathcal{T}_h)$, integrate over $K \in \mathcal{T}_h$ and use Green’s theorem. Summing over all $K \in \mathcal{T}_h$, we obtain the identity

$$\sum_{K \in \mathcal{T}_h} \int_K \nabla u \cdot \nabla v \, dx - \sum_{K \in \mathcal{T}_h} \int_{\partial K} (\mathbf{n}_K \cdot \nabla u) v \, dS = \int_{\Omega} f v \, dx, \quad (2.35)$$

where \mathbf{n}_K denotes the outer unit normal to ∂K . The surface integrals over ∂K make sense due to the regularity of u . (Since $u \in H^2(K)$, the derivatives $\partial u / \partial x_i$ have the trace on ∂K and $\partial u / \partial x_i|_{\partial K} \in L^2(\partial K)$ for $i = 1, \dots, d$; see Theorem 1.1 on traces.) We rewrite the surface integrals over ∂K according to the type of faces $\Gamma \in \mathcal{F}_h$ that form the boundary of the element $K \in \mathcal{T}_h$:

$$\begin{aligned} \sum_{K \in \mathcal{T}_h} \int_{\partial K} (\mathbf{n}_K \cdot \nabla u) v \, dS &= \sum_{\Gamma \in \mathcal{F}_h^D} \int_{\Gamma} (\mathbf{n}_{\Gamma} \cdot \nabla u) v \, dS + \sum_{\Gamma \in \mathcal{F}_h^N} \int_{\Gamma} (\mathbf{n}_{\Gamma} \cdot \nabla u) v \, dS \\ &\quad + \sum_{\Gamma \in \mathcal{F}_h^I} \int_{\Gamma} \mathbf{n}_{\Gamma} \cdot \left((\nabla u_{\Gamma}^{(L)}) v_{\Gamma}^{(L)} - (\nabla u_{\Gamma}^{(R)}) v_{\Gamma}^{(R)} \right) dS. \end{aligned} \quad (2.36)$$

(There is the sign “−” in the last integral, since \mathbf{n}_{Γ} is the outer unit normal to $\partial K_{\Gamma}^{(L)}$ but the inner unit normal to $\partial K_{\Gamma}^{(R)}$, see Sect. 2.3.1 or Fig. 2.2.)

Due to the assumption that $u \in H^2(\Omega)$, we have

$$[u]_{\Gamma} = [\nabla u]_{\Gamma} = 0, \quad \nabla u_{\Gamma}^{(L)} = \nabla u_{\Gamma}^{(R)} = \langle \nabla u \rangle_{\Gamma}, \quad \Gamma \in \mathcal{F}_h^I. \quad (2.37)$$

Thus, the integrand of the last integral in (2.36) can be written in the form

$$\mathbf{n}_\Gamma \cdot (\nabla u)_\Gamma^{(L)} v_\Gamma^{(L)} - \mathbf{n}_\Gamma \cdot (\nabla u)_\Gamma^{(R)} v_\Gamma^{(R)} = \mathbf{n}_\Gamma \cdot \langle \nabla u \rangle_\Gamma [v]_\Gamma. \quad (2.38)$$

By virtue of the Neumann boundary condition (2.1c),

$$\sum_{\Gamma \in \mathcal{F}_h^N} \int_\Gamma (\mathbf{n}_\Gamma \cdot \nabla u) v \, dS = \int_{\partial\Omega_N} g_N v \, dS. \quad (2.39)$$

Now, (2.33) and (2.35)–(2.39) imply that

$$\begin{aligned} & \sum_{K \in \mathcal{T}_h} \int_K \nabla u \cdot \nabla v \, dx - \sum_{\Gamma \in \mathcal{F}_h^I} \int_\Gamma \mathbf{n} \cdot \langle \nabla u \rangle [v] \, dS - \sum_{\Gamma \in \mathcal{F}_h^D} \int_\Gamma \mathbf{n} \cdot \nabla u \, v \, dS \\ &= \sum_{K \in \mathcal{T}_h} \int_K \nabla u \cdot \nabla v \, dx - \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_\Gamma \mathbf{n} \cdot \langle \nabla u \rangle [v] \, dS \\ &= \int_\Omega f v \, dx + \int_{\partial\Omega_N} g_N v \, dS, \quad v \in H^1(\Omega, \mathcal{T}_h). \end{aligned} \quad (2.40)$$

Here and in what follows, in integrals over Γ the symbol \mathbf{n} means \mathbf{n}_Γ .

Relation (2.40) is the basis of the DG discretization of problem (2.1). However, in order to guarantee the existence of the approximate solution and its convergence to the exact one, some additional terms have to be included in the DG formulation.

In order to mimic the continuity of the approximate solution in a weaker sense, we define the *interior and boundary penalty bilinear form*

$$\begin{aligned} J_h^\sigma(u, v) &= \sum_{\Gamma \in \mathcal{F}_h^I} \int_\Gamma \sigma[u][v] \, dS + \sum_{\Gamma \in \mathcal{F}_h^D} \int_\Gamma \sigma u v \, dS \\ &= \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_\Gamma \sigma[u][v] \, dS, \quad u, v \in H^1(\Omega, \mathcal{T}_h). \end{aligned} \quad (2.41)$$

The boundary penalty is associated with the boundary linear form

$$J_D^\sigma(v) = \sum_{\Gamma \in \mathcal{F}_h^D} \int_\Gamma \sigma u_D v \, dS. \quad (2.42)$$

Here $\sigma > 0$ is a penalty weight. Its choice will be discussed in Sect. 2.6. Obviously, for the exact strong solution $u \in H^2(\Omega)$,

$$J_h^\sigma(u, v) = J_D^\sigma(v) \quad \forall v \in H^1(\Omega, \mathcal{T}_h), \quad (2.43)$$

since $[u]_\Gamma = 0$ for $\Gamma \in \mathcal{F}_h^I$ and $[u]_\Gamma = u_\Gamma = u_D$ for $\Gamma \in \mathcal{F}_h^D$.

The interior penalty replaces the continuity of the approximate solution on interior faces, which is required in the standard conforming finite element method. The boundary penalty introduces the Dirichlet boundary condition in the discrete problem.

Moreover, the left-hand side of (2.40) is not symmetric with respect to u and v . In the theoretical analysis, it is advantageous to have some type of symmetry. Hence, it is desirable to include some additional term, which “symmetrizes” the left-hand side of (2.40) and which vanishes for the exact solution. Therefore, let $u \in H^1(\Omega) \cap H^2(\Omega, \mathcal{T}_h)$ be a function which satisfies the Dirichlet boundary condition (2.1b). Then we use the identity

$$\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \mathbf{n} \cdot \langle \nabla v \rangle [u] dS = \sum_{\Gamma \in \mathcal{F}_h^D} \int_{\Gamma} \mathbf{n} \cdot \nabla v u_D dS \quad \forall v \in H^2(\Omega, \mathcal{T}_h), \quad (2.44)$$

which is valid since $[u]_{\Gamma} = 0$ for $\Gamma \in \mathcal{F}_h^I$, $[u]_{\Gamma} = u_{\Gamma} = u_D$ for $\Gamma \in \mathcal{F}_h^D$ and $\langle \nabla v \rangle_{\Gamma} = \nabla v_{\Gamma}$ for $\Gamma \in \mathcal{F}_h^D$ by definition.

Now, without a deeper motivation, we introduce five variants of the *discontinuous Galerkin weak formulation*. Each particular method is commented on in Remark 2.10. Hence, we sum identity (2.40) with -1 , 1 or 0 -multiple of (2.44) and possibly add equality (2.43). This leads us to the following notation. For $u, v \in H^2(\Omega, \mathcal{T}_h)$ we introduce the bilinear *diffusion forms*

$$a_h^s(u, v) = \sum_{K \in \mathcal{T}_h} \int_K \nabla u \cdot \nabla v dx - \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} (\mathbf{n} \cdot \langle \nabla u \rangle [v] + \mathbf{n} \cdot \langle \nabla v \rangle [u]) dS, \quad (2.45a)$$

$$a_h^n(u, v) = \sum_{K \in \mathcal{T}_h} \int_K \nabla u \cdot \nabla v dx - \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} (\mathbf{n} \cdot \langle \nabla u \rangle [v] - \mathbf{n} \cdot \langle \nabla v \rangle [u]) dS, \quad (2.45b)$$

$$a_h^i(u, v) = \sum_{K \in \mathcal{T}_h} \int_K \nabla u \cdot \nabla v dx - \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \mathbf{n} \cdot \langle \nabla u \rangle [v] dS, \quad (2.45c)$$

and the right-hand side linear forms

$$F_h^s(v) = \int_{\Omega} f v dx + \sum_{\Gamma \in \mathcal{F}_h^N} \int_{\Gamma} g_N v dS - \sum_{\Gamma \in \mathcal{F}_h^D} \int_{\Gamma} \mathbf{n} \cdot \nabla v u_D dS, \quad (2.46a)$$

$$F_h^n(v) = \int_{\Omega} f v dx + \sum_{\Gamma \in \mathcal{F}_h^N} \int_{\Gamma} g_N v dS + \sum_{\Gamma \in \mathcal{F}_h^D} \int_{\Gamma} \mathbf{n} \cdot \nabla v u_D dS, \quad (2.46b)$$

$$F_h^i(v) = \int_{\Omega} f v dx + \sum_{\Gamma \in \mathcal{F}_h^N} \int_{\Gamma} g_N v dS. \quad (2.46c)$$

Moreover, for $u, v \in H^2(\Omega, \mathcal{T}_h)$ let us define the bilinear forms

$$A_h^s(u, v) = a_h^s(u, v), \quad (2.47a)$$

$$A_h^n(u, v) = a_h^n(u, v), \quad (2.47b)$$

$$A_h^{s,\sigma}(u, v) = a_h^s(u, v) + J_h^\sigma(u, v), \quad (2.47c)$$

$$A_h^{n,\sigma}(u, v) = a_h^n(u, v) + J_h^\sigma(u, v), \quad (2.47d)$$

$$A_h^{i,\sigma}(u, v) = a_h^i(u, v) + J_h^\sigma(u, v), \quad (2.47e)$$

and the linear forms

$$\ell_h^s(v) = F_h^s(v), \quad (2.48a)$$

$$\ell_h^n(v) = F_h^n(v), \quad (2.48b)$$

$$\ell_h^{s,\sigma}(v) = F_h^s(v) + J_D^\sigma(v), \quad (2.48c)$$

$$\ell_h^{n,\sigma}(v) = F_h^n(v) + J_D^\sigma(v), \quad (2.48d)$$

$$\ell_h^{i,\sigma}(v) = F_h^i(v) + J_D^\sigma(v). \quad (2.48e)$$

Since $S_{hp} \subset H^2(\Omega, \mathcal{T}_h)$, the forms (2.47) make sense for $u_h, v_h \in S_{hp}$. Consequently, we define five numerical schemes.

Definition 2.7 A function $u_h \in S_{hp}$ is called a *DG approximate solution* of problem (2.1), if it satisfies one of the following identities:

$$(i) \quad A_h^s(u_h, v_h) = \ell_h^s(v_h) \quad \forall v_h \in S_{hp}, \quad (2.49a)$$

$$(ii) \quad A_h^n(u_h, v_h) = \ell_h^n(v_h) \quad \forall v_h \in S_{hp}, \quad (2.49b)$$

$$(iii) \quad A_h^{s,\sigma}(u_h, v_h) = \ell_h^{s,\sigma}(v_h) \quad \forall v_h \in S_{hp}, \quad (2.49c)$$

$$(iv) \quad A_h^{n,\sigma}(u_h, v_h) = \ell_h^{n,\sigma}(v_h) \quad \forall v_h \in S_{hp}, \quad (2.49d)$$

$$(v) \quad A_h^{i,\sigma}(u_h, v_h) = \ell_h^{i,\sigma}(v_h) \quad \forall v_h \in S_{hp}, \quad (2.49e)$$

where the forms A_h^s, A_h^n, \dots , and $\ell_h^s, \ell_h^n, \dots$, are defined by (2.47) and (2.48), respectively.

The diffusion forms a_h^s, a_h^n, a_h^i defined by (2.45) can be simply written in the form

$$a_h(u, v) = \sum_{K \in \mathcal{T}_h} \int_K \nabla u \cdot \nabla v \, dx - \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} (\mathbf{n} \cdot \langle \nabla u \rangle [v] + \Theta \mathbf{n} \cdot \langle \nabla v \rangle [u]) \, dS, \quad (2.50)$$

where $\Theta = 1$ in the case of the form a_h^s , $\Theta = -1$ for a_h^n and $\Theta = 0$ for a_h^i and the bilinear forms $A_h^s, A_h^n, A_h^{s,\sigma}, A_h^{n,\sigma}$ and $A_h^{i,\sigma}$ defined by (2.47) can be written in the form

$$A_h(u, v) = a_h(u, v) + \vartheta J_h^\sigma(u, v), \quad (2.51)$$

where $\vartheta = 0$ for A_h^s and A_h^n and $\vartheta = 1$ for $A_h^{s,\sigma}$, $A_h^{n,\sigma}$ and $A_h^{i,\sigma}$.

Similarly we can write

$$F_h(v) = \int_{\Omega} f v \, dx + \sum_{\Gamma \in \mathcal{T}_h^N} \int_{\Gamma} g_N v \, dS - \Theta \sum_{\Gamma \in \mathcal{T}_h^D} \int_{\Gamma} \mathbf{n} \cdot \nabla v u_D \, dS, \quad (2.52)$$

with $\Theta = 1$ for F_h^s , $\Theta = -1$ for F_h^n and $\Theta = 0$ for F_h^i , and then the right-hand side form reads

$$\ell_h(v) = F_h(v) + \vartheta J_D^\sigma(v), \quad (2.53)$$

where $\vartheta = 0$ for ℓ_h^s and ℓ_h^n and $\vartheta = 1$ for $\ell_h^{s,\sigma}$, $\ell_h^{n,\sigma}$ and $\ell_h^{i,\sigma}$.

The form a_h^n ($\Theta = -1$), a_h^i ($\Theta = 0$) and a_h^s ($\Theta = 1$) represents the so-called *nonsymmetric*, *incomplete* and *symmetric* variant of the diffusion discretization, respectively.

If we denote by A_h any form defined by (2.47) and by ℓ_h , we denote the form defined by (2.53), i.e., any form given by (2.48), the *discrete problem* (2.49) can be formulated to find $u_h \in S_{hp}$ satisfying the identity

$$A_h(u_h, v_h) = \ell_h(v_h) \quad \forall v_h \in S_{hp}. \quad (2.54)$$

The discrete problem (2.54) is equivalent to a system of linear algebraic equations, which can be solved by a suitable direct or iterative method. Namely, let $\{\varphi_i, i = 1, \dots, N_h\}$ be a basis of the space S_{hp} , where $N_h = \dim S_{hp}$ (= dimension of S_{hp}). The approximate solution u_h is sought in the form $u_h(x) = \sum_{j=1}^{N_h} u^j \varphi_j(x)$, where u^j , $j = 1, \dots, N_h$, are unknown real coefficients. Then, due to the linearity of the form A_h , the discrete problem (2.54) is equivalent to the system

$$\sum_{j=1}^{N_h} A_h(\varphi_j, \varphi_i) u^j = \ell_h(\varphi_i), \quad i = 1, \dots, N_h. \quad (2.55)$$

It can be written in the matrix form

$$\mathbb{A}U = L,$$

where $\mathbb{A} = (a_{ij})_{i,j=1}^{N_h} = (A_h(\varphi_j, \varphi_i))_{i,j=1}^{N_h}$, $U = (u^j)_{j=1}^{N_h}$ and $L = (\ell_h(\varphi_j))_{j=1}^{N_h}$.

From the construction of the forms A_h and ℓ_h , one can see that the strong solution $u \in H^2(\Omega)$ of problem (2.1) satisfies the identity

$$A_h(u, v) = \ell_h(v) \quad \forall v \in H^2(\Omega, \mathcal{T}_h), \quad (2.56)$$

which represents the *consistency* of the method. Relations (2.54) and (2.56) imply the so-called *Galerkin orthogonality* of the error $e_h = u_h - u$ of the method:

$$A_h(e_h, v_h) = 0 \quad \forall v_h \in S_{hp}, \quad (2.57)$$

which will be used in analysing error estimates.

Remark 2.8 Comparing the above process of the derivation of the DG schemes with the abstract numerical method in Sect. 2.2, we see that we can define the function spaces

$$V = H^2(\Omega), \quad W_h = H^2(\Omega, \mathcal{T}_h), \quad V_h = S_{hp}. \quad (2.58)$$

However, as we will see later, the space W_h will not be equipped with the norm $\|\cdot\|_{H^2(\Omega, \mathcal{T}_h)}$ defined by (2.30), but by another norm introduced later in (2.103) will be used.

Remark 2.9 The interior and boundary penalty form J_h^σ together with the form J_D^σ replace the continuity of conforming finite element approximate solutions and represent Dirichlet boundary conditions. Thus, in contrast to standard conforming finite element techniques, both Dirichlet and Neumann boundary conditions are included automatically in the formulation (2.54) of the discrete problem. This is an advantage particularly in the case of nonhomogeneous Dirichlet boundary conditions, because it is not necessary to construct subsets of finite element spaces formed by functions approximating the Dirichlet boundary condition in a suitable way.

Remark 2.10 Method (2.49a) was introduced by Delves et al. ([76, 77, 172, 173]), who called it a *global element method*. Its advantage is the symmetry of the discrete problem due to the third term on the right-hand side of (2.45a). On the other hand, a significant disadvantage is that the bilinear form A_h^s is indefinite. This causes difficulties when dealing with time-dependent problems, because some eigenvalues of the operator associated with the form A_h can have negative real parts and then the resulting space-time discrete schemes become unconditionally unstable. Therefore, we prove in Lemma 2.36 the continuity of the bilinear form A_h^s , but further on we are not concerned with this method any more.

Scheme (2.49b) was introduced by Baumann and Oden in [12, 230] and is usually called the *Baumann–Oden method*. It is straightforward to show that the corresponding bilinear form A_h^n is positive semidefinite due to the third term on the right-hand side of (2.45b). An interesting property of this method is that it is unstable for piecewise linear approximations, i.e., for $p = 1$.

Scheme (2.49c) is called the *symmetric interior penalty Galerkin* (SIPG) method. It was derived by Arnold ([7]) and Wheeler ([283]) by adding penalty terms to the form A_h^s . (In this case a_h and F_h are defined by (2.50) and (2.52) with $\Theta = 1$.) This formulation leads to a symmetric bilinear form, which is coercive, if the penalty parameter σ is sufficiently large. Moreover, the Aubin–Nitsche duality technique

(also called Aubin–Nitsche trick) can be used to obtain an optimal error estimate in the $L^2(\Omega)$ -norm.

Method (2.49d), called the *nonsymmetric interior penalty Galerkin* (NIPG) method, was proposed by Girault, Rivi re and Wheeler in [239]. (Here $\Theta = -1$.) In this case the bilinear form $A_h^{n,\sigma}$ is nonsymmetric and does not allow one to obtain an optimal error estimate in the $L^2(\Omega)$ -norm with the aid of the Aubin–Nitsche trick. However, numerical experiments show that in some situations (for example, if uniform grids are used) the odd degrees of the polynomial approximation give the optimal order of convergence. On the other hand, a favorable property of the NIPG method is the coercivity of $A_h^{n,\sigma}(\cdot, \cdot)$ for any penalty parameter $\sigma > 0$.

Finally, method (2.49e), called the *incomplete interior penalty Galerkin* (IIPG) method ($\Theta = 0$), was studied in [74, 263, 265]. In this case the bilinear form $A_h^{i,\sigma}$ is nonsymmetric and does not allow one to obtain an optimal error estimate in the $L^2(\Omega)$ -norm. The penalty parameter σ has to be chosen sufficiently large in order to guarantee the coercivity of $A_h^{i,\sigma}$. The advantage of the IIPG method is the simplicity of the discrete diffusion operator, because the expressions from (2.44) do not appear in (2.45c). This is particularly advantageous in the case when the diffusion operator is nonlinear with respect to ∇u . (See, e.g., [87] or Chap. 9 of this book.)

It would also be possible to define the scheme $A_h^i(u, v) = \ell_h^i(v) \forall v \in S_{hp}$, where $A_h^i(u, v) = a_h^i(u, v)$ and $\ell_h^i(v) = F_h^i(v)$, but this method does not make sense, because it does not contain the Dirichlet boundary data u_D from condition (2.1b).

In the following, we deal with the theoretical analysis of the DGM applied to the numerical solution of the model problem (2.1). Namely, we pay attention to the existence and uniqueness of the approximate solution defined by (2.54) and derive error estimates.

2.5 Basic Tools of the Theoretical Analysis of DGM

Theoretical analysis of the DG method presented in this book is based on three fundamental tools: the *multiplicative trace inequality*, the *inverse inequality*, and the *approximation properties* of the spaces of piecewise polynomial functions. In this section we introduce and prove these important tools under the assumptions about the meshes in Sect. 2.3.2.

Our first objective will be to summarize some important concepts and results from finite element theory, treated, e.g., in [52].

Definition 2.11 Let $n > 0$ be an integer. We say that sets $\omega, \widehat{\omega} \subset \mathbb{R}^n$ are *affine equivalent*, if there exists an invertible affine mapping $F_\omega : \widehat{\omega} \rightarrow \omega$ such that $F_\omega(\widehat{\omega}) = \omega$ and

$$x = F_\omega(\hat{x}) = \mathbb{B}_\omega \hat{x} + b_\omega \in \omega, \quad \hat{x} \in \widehat{\omega}, \quad (2.59)$$

where \mathbb{B}_ω is an $n \times n$ nonsingular matrix and $b_\omega \in \mathbb{R}^n$.

If $\hat{v} : \hat{\omega} \rightarrow \mathbb{R}$, then the inverse mapping F_ω^{-1} allows us to transform the function \hat{v} to $v : \omega \rightarrow \mathbb{R}$ by the relation

$$v(x) = \hat{v}(F_\omega^{-1}(x)), \quad x \in \omega. \quad (2.60)$$

Hence,

$$v = \hat{v} \circ F_\omega^{-1}, \quad \hat{v} = v \circ F_\omega \quad (2.61)$$

and

$$\hat{v}(\hat{x}) = v(x) \text{ for all } \hat{x}, x \text{ in the correspondence (2.59).}$$

If \mathbb{B} is an $n \times n$ matrix, then its norm associated with the Euclidean norm $|\cdot|$ in \mathbb{R}^n is defined as $\|\mathbb{B}\| = \sup_{0 \neq x \in \mathbb{R}^n} |\mathbb{B}x|/|x|$.

The following lemmas give us bounds for the norms of matrices \mathbb{B}_ω and \mathbb{B}_ω^{-1} and the relations between Sobolev seminorms of functions v and \hat{v} satisfying (2.61). First, we introduce the following notation for bounded domains $\omega, \hat{\omega}$:

$$h_\omega = \text{diam}(\omega), \quad h_{\hat{\omega}} = \text{diam}(\hat{\omega}), \quad (2.62)$$

$$\rho_\omega = \text{radius of the largest ball inscribed into } \overline{\omega}, \quad (2.63)$$

$$\rho_{\hat{\omega}} = \text{radius of the largest ball inscribed into } \overline{\hat{\omega}}.$$

Lemma 2.12 *Let $\omega, \hat{\omega} \subset \mathbb{R}^n$ be affine-equivalent bounded domains with the invertible mapping $F_\omega(\hat{x}) = \mathbb{B}_\omega \hat{x} + b_\omega \in \omega$ for $\hat{x} \in \hat{\omega}$. Then*

$$\|\mathbb{B}_\omega\| \leq \frac{h_\omega}{2\rho_{\hat{\omega}}}, \quad \|\mathbb{B}_\omega^{-1}\| \leq \frac{h_{\hat{\omega}}}{2\rho_\omega}. \quad (2.64)$$

Further, the substitution theorem implies that

$$|\det(\mathbb{B}_\omega)| = |\omega|/|\hat{\omega}|, \quad (2.65)$$

where $|\omega|$ and $|\hat{\omega}|$ denote the n -dimensional Lebesgue measure of ω and $\hat{\omega}$, respectively.

For the proof of (2.64) see [52, Theorem 3.1.3]. The proof of (2.65) is a consequence of the substitution theorem. Further, we cite here Theorem 3.1.2 from [52].

Lemma 2.13 *Let $\omega, \hat{\omega} \subset \mathbb{R}^n$ be affine-equivalent bounded domains with the invertible mapping $F_\omega(\hat{x}) = \mathbb{B}_\omega \hat{x} + b_\omega \in \omega$ for $\hat{x} \in \hat{\omega}$. If $v \in W^{m,\alpha}(\omega)$ for some integer $m \geq 0$ and some $\alpha \in [1, \infty]$, then the function $\hat{v} = v \circ F_\omega \in W^{m,\alpha}(\hat{\omega})$. Moreover, there exists a constant C depending on m and d only such that*

$$|\hat{v}|_{W^{m,\alpha}(\hat{\omega})} \leq C \|\mathbb{B}_\omega\|^m |\det(\mathbb{B}_\omega)|^{-1/\alpha} |v|_{W^{m,\alpha}(\omega)}, \quad (2.66)$$

$$|v|_{W^{m,\alpha}(\omega)} \leq C \|\mathbb{B}_\omega^{-1}\|^m |\det(\mathbb{B}_\omega)|^{1/\alpha} |\hat{v}|_{W^{m,\alpha}(\hat{\omega})}. \quad (2.67)$$

In our finite element analysis, we have $n = d$ and the set ω represents an element $K \in \mathcal{T}_h$ and $\widehat{\omega}$ is chosen as a reference element \widehat{K} , i.e., the simplex with vertices

$$\begin{aligned} \hat{a}_1 = (0, 0, \dots, 0), \quad \hat{a}_2 = (1, 0, \dots, 0), \quad \hat{a}_3 = (0, 1, 0, \dots, 0), \dots \\ \dots, \quad \hat{a}_{d+1} = (0, 0, \dots, 1) \in \mathbb{R}^d. \end{aligned} \quad (2.68)$$

The elements K and \widehat{K} are considered as closed sets. The Sobolev spaces over K and \widehat{K} are defined as the spaces over the interiors of these sets. (In Sect. 7.3, we will also apply the above results to the case with $n = 1$, $\omega = \Gamma \in \mathcal{F}_h$ and $\widehat{\omega} = (0, 1)$.)

As a consequence of the above results we can formulate the following assertions.

Corollary 2.14 *If $K \in \mathcal{T}_h$ and $v \in H^m(K)$, where $m \geq 0$ is an integer, then the function $\hat{v}(\hat{x}) = v(F_K(\hat{x})) \in H^m(\widehat{K})$ and*

$$|v|_{H^m(K)} \leq c_c h_K^{\frac{d}{2}-m} |\hat{v}|_{H^m(\widehat{K})}, \quad (2.69)$$

$$|\hat{v}|_{H^m(\widehat{K})} \leq c_c h_K^{m-\frac{d}{2}} |v|_{H^m(K)}, \quad (2.70)$$

where $c_c > 0$ depends on the shape regularity constant C_R but not on K and v .

Exercises 2.15 Prove (2.69) and (2.70) using the shape-regularity assumption (2.19) and the results of Lemmas 2.12 and 2.13.

In deriving error estimates we apply the following important result from [52, Theorem 3.1.4].

Theorem 2.16 *Let $\widehat{\omega} \subset \mathbb{R}^n$ be a bounded domain and for some integers $p \geq 0$ and $m \geq 0$ and some numbers $\alpha, \beta \in [1, \infty]$, let the spaces $W^{p+1,\alpha}(\widehat{\omega})$ and $W^{m,\beta}(\widehat{\omega})$ satisfy the continuous embedding*

$$W^{p+1,\alpha}(\widehat{\omega}) \hookrightarrow W^{m,\beta}(\widehat{\omega}). \quad (2.71)$$

Let $\widehat{\Pi}$ be a continuous linear mapping of $W^{p+1,\alpha}(\widehat{\omega})$ into $W^{m,\beta}(\widehat{\omega})$ such that

$$\widehat{\Pi}\hat{\phi} = \hat{\phi} \quad \forall \hat{\phi} \in P_p(\widehat{\omega}). \quad (2.72)$$

Let a set ω be affine-equivalent to the set $\widehat{\omega}$. This means that there exists an affine mapping $x = F_\omega(\hat{x}) = \mathbb{B}_\omega \hat{x} + b_\omega \in \omega$ for $\hat{x} \in \widehat{\omega}$, where \mathbb{B}_ω is a nonsingular $n \times n$ matrix and $b_\omega \in \mathbb{R}^n$. Let the mapping Π_ω be defined by

$$\Pi_\omega v(x) = (\widehat{\Pi}\hat{v})(F_\omega^{-1}(x)), \quad (2.73)$$

for all functions $\hat{v} \in W^{p+1,\alpha}(\hat{\omega})$ and $v \in W^{p+1,\alpha}(\omega)$ such that $\hat{v}(\hat{x}) = v(F_\omega(\hat{x})) = v(x)$. Then there exists a constant $C(\hat{\Pi}, \hat{\omega})$ such that

$$|\hat{\Pi}\hat{v} - \hat{v}|_{W^{m,\beta}(\hat{\omega})} \leq C(\hat{\Pi}, \hat{\omega}) |\hat{v}|_{W^{p+1,\alpha}(\hat{\omega})}, \quad (2.74)$$

and

$$|v - \Pi_\omega v|_{W^{m,\beta}(\omega)} \leq C(\hat{\Pi}, \hat{\omega}) |\omega|^{(1/\beta)-(1/\alpha)} \frac{h_\omega^{p+1}}{\rho_\omega^m} |v|_{W^{p+1,\alpha}(\omega)} \quad (2.75)$$

$$\forall v \in W^{p+1,\alpha}(\omega),$$

with $h_\omega = \text{diam}(\omega)$, ρ_ω defined as the radius of the largest ball inscribed into $\bar{\omega}$ and $|\omega|$ defined as the n -dimensional Lebesgue measure of the set ω . We set $1/\infty := 0$.

Exercises 2.17 Prove (2.75) using (2.74), (2.66), (2.67), (2.64) and (2.65).

Another important result used often in finite element theory is the Bramble–Hilbert lemma (see [52, Theorem 4.1.3] or [287, Theorem 9.3]).

Theorem 2.18 (Bramble–Hilbert lemma) *Let us assume that $\omega \subset \mathbb{R}^n$ is a bounded domain with Lipschitz boundary. Let $p \geq 0$ be an integer and $\alpha \in [1, \infty]$ and let f be a continuous linear functional on the space $W^{p+1,\alpha}(\Omega)$ (i.e., $f \in (W^{p+1,\alpha}(\omega))^*$) satisfying the condition*

$$f(v) = 0 \quad \forall v \in P_p(\omega). \quad (2.76)$$

Then there exists a constant $C_{BH} > 0$ depending only on ω such that

$$|f(v)| \leq C_{BH} \|f\|_{(W^{p+1,\alpha}(\omega))^*} |v|_{W^{p+1,\alpha}(\omega)} \quad \forall v \in W^{p+1,\alpha}(\omega). \quad (2.77)$$

2.5.1 Multiplicative Trace Inequality

The forms a_h and J_h^σ given by (2.45) and (2.41), respectively, contain several integrals over faces. Therefore, in the theoretical analysis we need to estimate norms over faces by norms over elements. These estimates are usually obtained using the *multiplicative trace inequality*. In the literature, it is possible to find several variants of the multiplicative trace inequality. Here, we present the variant, which suits our considerations.

Lemma 2.19 (Multiplicative trace inequality) *Let the shape-regularity assumption (2.19) be satisfied. Then there exists a constant $C_M > 0$ independent of v , h and K such that*

$$\|v\|_{L^2(\partial K)}^2 \leq C_M \left(\|v\|_{L^2(K)} \|v\|_{H^1(K)} + h_K^{-1} \|v\|_{L^2(K)}^2 \right), \quad (2.78)$$

$$K \in \mathcal{T}_h, \quad v \in H^1(K), \quad h \in (0, \tilde{h}).$$

Proof Let $K \in \mathcal{T}_h$ be arbitrary but fixed. We denote by \mathbf{x}_K the center of the largest d -dimensional ball inscribed into the simplex K . Without loss of generality we suppose that \mathbf{x}_K is the origin of the coordinate system.

Since the space $C^\infty(K)$ is dense in $H^1(K)$, it is sufficient to prove (2.78) for $v \in C^\infty(K)$. We start from the following relation obtained from Green's identity (1.23):

$$\int_{\partial K} v^2 \mathbf{x} \cdot \mathbf{n} \, dS = \int_K \nabla \cdot (v^2 \mathbf{x}) \, dx, \quad v \in C^\infty(K), \quad (2.79)$$

where \mathbf{n} denotes here the outer unit normal to ∂K . Let \mathbf{n}_Γ be the outer unit normal to K on a side Γ of K . Then

$$\mathbf{x} \cdot \mathbf{n}_\Gamma = |\mathbf{x}| |\mathbf{n}_\Gamma| \cos \alpha = |\mathbf{x}| \cos \alpha = \rho_K, \quad \mathbf{x} \in \Gamma, \quad (2.80)$$

see Fig. 2.3. From (2.80) we have

$$\int_{\partial K} v^2 \mathbf{x} \cdot \mathbf{n} \, dS = \sum_{\Gamma \subset \partial K} \int_{\Gamma} v^2 \mathbf{x} \cdot \mathbf{n}_\Gamma \, dS = \rho_K \sum_{\Gamma \subset \partial K} \int_{\Gamma} v^2 \, dS = \rho_K \|v\|_{L^2(\partial K)}^2. \quad (2.81)$$

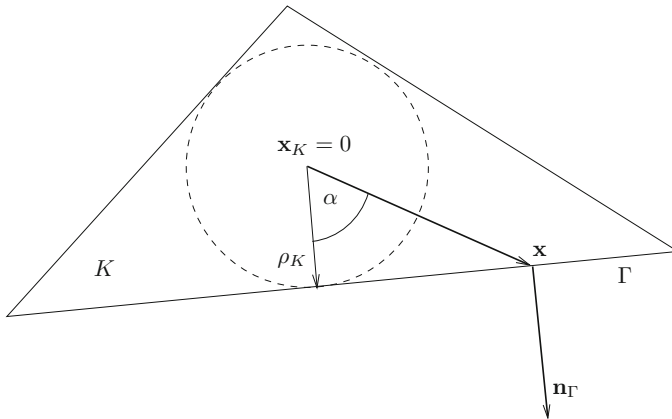


Fig. 2.3 Simplex K with its face Γ

Moreover,

$$\begin{aligned} \int_K \nabla \cdot (v^2 \mathbf{x}) \, dx &= \int_K \left(v^2 \nabla \cdot \mathbf{x} + \mathbf{x} \cdot \nabla v^2 \right) \, dx \\ &= d \int_K v^2 \, dx + 2 \int_K v \mathbf{x} \cdot \nabla v \, dx \leq d \|v\|_{L^2(K)}^2 + 2 \int_K |v \mathbf{x} \cdot \nabla v| \, dx. \end{aligned} \quad (2.82)$$

With the aid of the Cauchy inequality, the second term of (2.82) is estimated as

$$2 \int_K |v \mathbf{x} \cdot \nabla v| \, dx \leq 2 \sup_{\mathbf{x} \in K} |\mathbf{x}| \int_K |v| |\nabla v| \, dx \leq 2 h_K \|v\|_{L^2(K)} |v|_{H^1(K)}. \quad (2.83)$$

Then (2.19), (2.79), (2.81)–(2.83) give

$$\begin{aligned} \|v\|_{L^2(\partial K)}^2 &\leq \frac{1}{\rho_K} \left[2 h_K \|v\|_{L^2(K)} |v|_{H^1(K)} + d \|v\|_{L^2(K)}^2 \right] \\ &\leq C_R \left[2 \|v\|_{L^2(K)} |v|_{H^1(K)} + \frac{d}{h_K} \|v\|_{L^2(K)}^2 \right], \end{aligned} \quad (2.84)$$

which proves (2.78) with $C_M = C_R \max\{2, d\}$. \square

Exercises 2.20 Prove that the multiplicative trace inequality is valid also for vector-valued functions $\mathbf{v} : \Omega \rightarrow \mathbb{R}^n$, i.e.,

$$\|\mathbf{v}\|_{L^2(\partial K)}^2 \leq C_M \left(\|\mathbf{v}\|_{L^2(K)} |\mathbf{v}|_{H^1(K)} + h_K^{-1} \|\mathbf{v}\|_{L^2(K)}^2 \right), \quad \mathbf{v} \in (H^1(K))^n, \quad K \in \mathcal{T}_h. \quad (2.85)$$

Hint: Use (2.78) for each component of $\mathbf{v} = (v_1, \dots, v_n)$, sum these inequalities and apply the discrete Cauchy inequality (1.52).

2.5.2 Inverse Inequality

In deriving error estimates, we need to estimate the H^1 -seminorm of a polynomial function by its L^2 -norm, i.e., we apply the so-called *inverse inequality*.

Lemma 2.21 (Inverse inequality) *Let the shape-regularity assumption (2.19) be satisfied. Then there exists a constant $C_I > 0$ independent of v , h and K such that*

$$|v|_{H^1(K)} \leq C_I h_K^{-1} \|v\|_{L^2(K)} \quad \forall v \in P_p(K), \quad \forall K \in \mathcal{T}_h, \quad \forall h \in (0, \bar{h}). \quad (2.86)$$

Proof Let \widehat{K} be a reference triangle and $F_K : \widehat{K} \rightarrow K$, $K \in \mathcal{T}_h$ be an affine mapping such that $F_K(\widehat{K}) = K$. By (2.69) (for $m = 1$) and (2.70) (for $m = 0$) we have

$$|v|_{H^1(K)} \leq c_c h_K^{\frac{d}{2}-1} |\hat{v}|_{H^1(\hat{K})}, \quad \|\hat{v}\|_{L^2(\hat{K})} \leq c_c h_K^{-\frac{d}{2}} \|v\|_{L^2(K)}. \quad (2.87)$$

From [253, Theorem 4.76], we have

$$|\hat{v}|_{H^1(\hat{K})} \leq c_s p^2 \|\hat{v}\|_{L^2(\hat{K})}, \quad \hat{v} \in P_p(\hat{K}), \quad (2.88)$$

where $c_s > 0$ depends on d but not on \hat{v} and p . A simple combination of (2.87) and (2.88) proves (2.86) with $C_I = c_s c_c^2 p^2$. Let us note that (2.88) is a consequence of the norm equivalence on finite-dimensional spaces. \square

Other inverse inequalities will appear in Sect. 7.3, Lemma 7.35.

2.5.3 Approximation Properties

With respect to the error analysis of the abstract numerical method treated in Sect. 2.2, a suitable S_{hp} -interpolation has to be introduced. Let \mathcal{T}_h be a given triangulation of the domain Ω . Then for each $K \in \mathcal{T}_h$, we define the mapping $\pi_{K,p} : L^2(K) \rightarrow P_p(K)$ such that for every $\varphi \in L^2(K)$

$$\pi_{K,p} \varphi \in P_p(K), \quad \int_K (\pi_{K,p} \varphi) v \, dx = \int_K \varphi v \, dx \quad \forall v \in P_p(K). \quad (2.89)$$

On the basis of the mappings $\pi_{K,p}$ we introduce the S_{hp} -interpolation Π_{hp} , defined for all $\varphi \in L^2(\Omega)$ by

$$(\Pi_{hp} \varphi)|_K = \pi_{K,p}(\varphi|_K) \quad \forall K \in \mathcal{T}_h. \quad (2.90)$$

It can be easily shown that if $\varphi \in L^2(\Omega)$, then

$$\Pi_{hp} \varphi \in S_{hp}, \quad \int_{\Omega} (\Pi_{hp} \varphi) v \, dx = \int_{\Omega} \varphi v \, dx \quad \forall v \in S_{hp}. \quad (2.91)$$

Hence, Π_{hp} is the $L^2(\Omega)$ -projection on the space S_{hp} .

The approximation properties of the interpolation operators $\pi_{K,p}$ and Π_{hp} are the consequence of Theorem 2.16.

Lemma 2.22 *Let the shape-regularity assumption (2.19) be valid and let p, q, s be integers, $p \geq 0$, $0 \leq q \leq \mu$, where $\mu = \min(p+1, s)$. Then there exists a constant $C_A > 0$ such that*

$$|\pi_{K,p} v - v|_{H^q(K)} \leq C_A h_K^{\mu-q} |v|_{H^\mu(K)} \quad \forall v \in H^s(K) \quad \forall K \in \mathcal{T}_h \quad \forall h \in (0, \bar{h}). \quad (2.92)$$

Hence, if $p \geq 1$ and $s \geq 2$, then

$$\|\pi_{K,p}v - v\|_{L^2(K)} \leq C_A h_K^\mu |v|_{H^\mu(K)} \quad \forall v \in H^s(K) \quad \forall K \in \mathcal{T}_h \quad \forall h \in (0, \bar{h}), \quad (2.93)$$

$$|\pi_{K,p}v - v|_{H^1(K)} \leq C_A h_K^{\mu-1} |v|_{H^\mu(K)} \quad \forall v \in H^s(K) \quad \forall K \in \mathcal{T}_h \quad \forall h \in (0, \bar{h}), \quad (2.94)$$

$$|\pi_{K,p}v - v|_{H^2(K)} \leq C_A h_K^{\mu-2} |v|_{H^\mu(K)} \quad \forall v \in H^s(K) \quad \forall K \in \mathcal{T}_h \quad \forall h \in (0, \bar{h}). \quad (2.95)$$

Moreover, we have

$$\|\pi_{K,1}v - v\|_{L^\infty(K)} \leq C_A h_K |v|_{W^{1,\infty}(K)} \quad \forall v \in W^{1,\infty}(K) \quad \forall K \in \mathcal{T}_h \quad \forall h \in (0, \bar{h}). \quad (2.96)$$

Exercises 2.23 Prove Lemma 2.22 using Theorem 2.16 and assumption (2.19).

The above results immediately imply the approximation properties of the operator Π_{hp} .

Lemma 2.24 *Let the shape-regularity assumption (2.19) be satisfied and let p, q, s be integers, $p \geq 0$, $0 \leq q \leq \mu$, where $\mu = \min(p+1, s)$. Then*

$$|\Pi_{hp}v - v|_{H^q(\Omega, \mathcal{T}_h)} \leq C_A h^{\mu-q} |v|_{H^\mu(\Omega, \mathcal{T}_h)}, \quad v \in H^s(\Omega, \mathcal{T}_h), \quad h \in (0, \bar{h}), \quad (2.97)$$

where C_A is the constant from (2.92). Hence, if $p \geq 1$ and $s \geq 2$, then

$$\|\Pi_{hp}v - v\|_{L^2(\Omega)} \leq C_A h^\mu |v|_{H^\mu(\Omega, \mathcal{T}_h)}, \quad v \in H^s(\Omega, \mathcal{T}_h), \quad h \in (0, \bar{h}), \quad (2.98)$$

$$|\Pi_{hp}v - v|_{H^1(\Omega, \mathcal{T}_h)} \leq C_A h^{\mu-1} |v|_{H^\mu(\Omega, \mathcal{T}_h)}, \quad v \in H^s(\Omega, \mathcal{T}_h), \quad h \in (0, \bar{h}), \quad (2.99)$$

$$|\Pi_{hp}v - v|_{H^2(\Omega, \mathcal{T}_h)} \leq C_A h^{\mu-2} |v|_{H^\mu(\Omega, \mathcal{T}_h)}, \quad v \in H^s(\Omega, \mathcal{T}_h), \quad h \in (0, \bar{h}). \quad (2.100)$$

Proof Using (2.90), definition of the seminorm in a broken Sobolev space (2.31) and the approximation properties (2.92), we obtain (2.97). This immediately implies (2.98)–(2.100). \square

Moreover, using the combination of the multiplicative trace inequality (2.78) and Lemma 2.22, we can prove the approximation properties of the operator Π_{hp} in the norms defined over the boundaries of elements.

Lemma 2.25 *Let the shape-regularity assumption (2.19) be satisfied and let $p \geq 1$, $s \geq 2$ be integers and $\alpha \geq -1$. Then*

$$\sum_{K \in \mathcal{T}_h} h_K^\alpha \|\Pi_{hp} v - v\|_{L^2(\partial K)}^2 \leq 2C_M C_A^2 h^{2\mu-1+\alpha} |v|_{H^\mu(\Omega, \mathcal{T}_h)}^2, \quad (2.101)$$

$$\sum_{K \in \mathcal{T}_h} h_K^\alpha \|\nabla(\Pi_{hp} v - v)\|_{L^2(\partial K)}^2 \leq 2C_M C_A^2 h^{2\mu-3+\alpha} |v|_{H^\mu(\Omega, \mathcal{T}_h)}^2, \quad (2.102)$$

$$v \in H^s(\Omega, \mathcal{T}_h), \quad h \in (0, \bar{h}),$$

where $\mu = \min(p+1, s)$, C_M is the constant from (2.78) and C_A is the constant from (2.92).

Proof (i) Let $v \in H^s(\Omega, \mathcal{T}_h)$. For simplicity we put $\eta = \Pi_{hp} v - v$. Then relation (2.90) implies that $\eta|_K = \pi_{K,p} v|_K - v|_K$ for $K \in \mathcal{T}_h$. Using the multiplicative trace inequality (2.78), the approximation property (2.92), and the seminorm definition (2.31), we have

$$\begin{aligned} \sum_{K \in \mathcal{T}_h} h_K^\alpha \|\eta\|_{L^2(\partial K)}^2 &\leq C_M \sum_{K \in \mathcal{T}_h} h_K^\alpha \left(\|\eta\|_{L^2(K)} \|\eta\|_{H^1(K)} + h_K^{-1} \|\eta\|_{L^2(K)}^2 \right) \\ &\leq C_M \sum_{K \in \mathcal{T}_h} h_K^\alpha C_A^2 \left(h_K^\mu h_K^{\mu-1} + h_K^{-1} h_K^{2\mu} \right) |v|_{H^\mu(K)}^2 \\ &\leq 2C_M C_A^2 h^{2\mu-1+\alpha} |v|_{H^\mu(\Omega, \mathcal{T}_h)}^2. \end{aligned}$$

(ii) Similarly as above, using the vector-valued variant of the multiplicative trace inequality (2.85), identities (1.21) and the approximation property (2.92), we get

$$\begin{aligned} \sum_{K \in \mathcal{T}_h} h_K^\alpha \|\nabla \eta\|_{L^2(\partial K)}^2 &\leq C_M \sum_{K \in \mathcal{T}_h} h_K^\alpha \left(\|\nabla \eta\|_{L^2(K)} \|\nabla \eta\|_{H^1(K)} + h_K^{-1} \|\nabla \eta\|_{L^2(K)}^2 \right) \\ &= C_M \sum_{K \in \mathcal{T}_h} h_K^\alpha \left(\|\eta\|_{H^1(K)} \|\eta\|_{H^2(K)} + h_K^{-1} \|\eta\|_{H^1(K)}^2 \right) \\ &\leq C_M \sum_{K \in \mathcal{T}_h} h_K^\alpha C_A^2 \left(h_K^{\mu-1} h_K^{\mu-2} + h_K^{-1} h_K^{2(\mu-1)} \right) |v|_{H^\mu(K)}^2 \\ &\leq 2C_M C_A^2 h^{2\mu-3+\alpha} |v|_{H^\mu(\Omega, \mathcal{T}_h)}^2. \end{aligned}$$

□

2.6 Existence and Uniqueness of the Approximate Solution

We start with the theoretical analysis of the DGM, namely we prove the existence of a numerical solution defined by (2.54). Then, in Sect. 2.7, we derive error estimates. We follow the formal analysis of the abstract numerical methods in Sect. 2.2. Therefore, we show the *continuity* and the *coercivity* of the form A_h given by (2.47) in a suitable

norm. This norm should reflect the discontinuity of functions from the broken Sobolev spaces $H^1(\Omega, \mathcal{T}_h)$. To this end, we define the following mesh-dependent norm

$$\|u\|_{\mathcal{T}_h} = \left(|u|_{H^1(\Omega, \mathcal{T}_h)}^2 + J_h^\sigma(u, u) \right)^{1/2}, \quad (2.103)$$

where $|\cdot|_{H^1(\Omega, \mathcal{T}_h)}$ and J_h^σ are given by (2.31) and (2.41), respectively.

In what follows, because there is no danger of misunderstanding, we omit the subscript \mathcal{T}_h . This means that we simply write $\|\cdot\| = \|\cdot\|_{\mathcal{T}_h}$. We call $\|\cdot\|$ the *DG-norm*.

Exercises 2.26 Prove that $\|\cdot\|$ is a norm in the spaces $H^1(\Omega, \mathcal{T}_h)$ and S_{hp} .

2.6.1 The Choice of Penalty Weight σ

In the following considerations we assume that the system $\{\mathcal{T}_h\}_{h \in (0, \bar{h})}$ of triangulations satisfies the shape-regularity assumption (2.19) and the equivalence condition (2.20).

We consider the penalty weight $\sigma : \cup_{\Gamma \in \mathcal{F}_h^{ID}} \rightarrow \mathbb{R}$ in the form

$$\sigma|_\Gamma = \sigma_\Gamma = \frac{C_W}{h_\Gamma}, \quad \Gamma \in \mathcal{F}_h^{ID}, \quad (2.104)$$

where $C_W > 0$ is the *penalization constant* and $h_\Gamma (\sim h)$ is the quantity given by one of the possibilities from (2.24)–(2.27) with respect to the considered mesh assumptions (MA1)–(MA4), see Lemma 2.5. Let us note that in some cases it is possible to consider a different form of the penalty parameter σ , as mentioned in Remark 2.51.

Under the introduced notation, in view of (2.41), (2.42) and (2.104), the interior and boundary penalty form and the associated boundary linear form read as

$$J_h^\sigma(u, v) = \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_\Gamma \frac{C_W}{h_\Gamma} [u] [v] dS, \quad J_D^\sigma(v) = \sum_{\Gamma \in \mathcal{F}_h^D} \int_\Gamma \frac{C_W}{h_\Gamma} u_D v dS. \quad (2.105)$$

In what follows, we introduce technical lemmas, which will be useful in the theoretical analysis.

Lemma 2.27 *Let (2.20) be valid. Then for each $v \in H^1(\Omega, \mathcal{T}_h)$ we have*

$$\sum_{\Gamma \in \mathcal{F}_h^{ID}} h_\Gamma^{-1} \int_\Gamma [v]^2 dS \leq \frac{2}{C_T} \sum_{K \in \mathcal{T}_h} h_K^{-1} \int_{\partial K} |v|^2 dS, \quad (2.106)$$

$$\sum_{\Gamma \in \mathcal{F}_h^{ID}} h_\Gamma \int_\Gamma \langle v \rangle^2 dS \leq C_G \sum_{K \in \mathcal{T}_h} h_K \int_{\partial K} |v|^2 dS. \quad (2.107)$$

Hence,

$$\sum_{\Gamma \in \mathcal{F}_h^{ID}} \sigma_\Gamma \| [v] \|_{L^2(\Gamma)}^2 \leq \frac{2C_W}{C_T} \sum_{K \in \mathcal{T}_h} h_K^{-1} \| v \|_{L^2(\partial K)}^2, \quad (2.108)$$

$$\sum_{\Gamma \in \mathcal{F}_h^{ID}} \frac{1}{\sigma_\Gamma} \| \langle v \rangle \|_{L^2(\Gamma)}^2 \leq \frac{C_G}{C_W} \sum_{K \in \mathcal{T}_h} h_K \| v \|_{L^2(\partial K)}^2. \quad (2.109)$$

Proof (i) By definition (2.32), the inequality

$$(\gamma + \delta)^2 \leq 2(\gamma^2 + \delta^2), \quad \gamma, \delta \in \mathbb{R}, \quad (2.110)$$

and (2.20) we have

$$\begin{aligned} & \sum_{\Gamma \in \mathcal{F}_h^{ID}} h_\Gamma^{-1} \int_\Gamma [v]^2 dS \\ &= \sum_{\Gamma \in \mathcal{F}_h^I} h_\Gamma^{-1} \int_\Gamma \left| v_\Gamma^{(L)} - v_\Gamma^{(R)} \right|^2 dS + \sum_{\Gamma \in \mathcal{F}_h^D} h_\Gamma^{-1} \int_\Gamma \left| v_\Gamma^{(L)} \right|^2 dS \\ &\leq 2 \sum_{\Gamma \in \mathcal{F}_h^I} h_\Gamma^{-1} \int_\Gamma \left(\left| v_\Gamma^{(L)} \right|^2 + \left| v_\Gamma^{(R)} \right|^2 \right) dS + \sum_{\Gamma \in \mathcal{F}_h^D} h_\Gamma^{-1} \int_\Gamma \left| v_\Gamma^{(L)} \right|^2 dS \\ &\leq 2C_T^{-1} \sum_{\Gamma \in \mathcal{F}_h^{ID}} h_{K_\Gamma}^{-1} \int_\Gamma \left| v_\Gamma^{(L)} \right|^2 dS + 2C_T^{-1} \sum_{\Gamma \in \mathcal{F}_h^I} h_{K_\Gamma}^{-1} \int_\Gamma \left| v_\Gamma^{(R)} \right|^2 dS \\ &\leq 2C_T^{-1} \sum_{K \in \mathcal{T}_h} h_K^{-1} \int_{\partial K} |v|^2 dS. \end{aligned}$$

This and (2.104) immediately imply (2.108).

(ii) In the proof of (2.107) we proceed similarly, using (2.32), (2.20) and (2.110). Inequalities (2.108) and (2.109) are obtained from (2.106), (2.107) and (2.104). \square

2.6.2 Continuity of Diffusion Bilinear Forms

First, we prove several auxiliary assertions.

Lemma 2.28 Any form a_h defined by (2.45) satisfies the estimate

$$|a_h(u, v)| \leq \|u\|_{1,\sigma} \|v\|_{1,\sigma} \quad \forall u, v \in H^2(\Omega, \mathcal{T}_h), \quad (2.111)$$

where

$$\begin{aligned} \|v\|_{1,\sigma}^2 &= \|v\|^2 + \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla v \rangle)^2 dS \\ &= |v|_{H^1(\Omega, \mathcal{T}_h)}^2 + J_h^\sigma(v, v) + \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla v \rangle)^2 dS. \end{aligned} \quad (2.112)$$

Proof It follows from (2.45) that

$$\begin{aligned} |a_h(u, v)| &\leq \underbrace{\sum_{K \in \mathcal{T}_h} \int_K |\nabla u \cdot \nabla v| dx}_{\chi_1} \\ &\quad + \underbrace{\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} |\mathbf{n} \cdot \langle \nabla u \rangle [v]| dS}_{\chi_2} + \underbrace{\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} |\mathbf{n} \cdot \langle \nabla v \rangle [u]| dS}_{\chi_3}. \end{aligned} \quad (2.113)$$

(For the form a_h^i the term χ_3 vanishes, of course.) Obviously, the Cauchy inequality, the discrete Cauchy inequality, and (2.31) imply that

$$\chi_1 \leq \sum_{K \in \mathcal{T}_h} |u|_{H^1(K)} |v|_{H^1(K)} \leq |u|_{H^1(\Omega, \mathcal{T}_h)} |v|_{H^1(\Omega, \mathcal{T}_h)}. \quad (2.114)$$

Further, by the Cauchy inequality,

$$\begin{aligned} \chi_2 &\leq \sum_{\Gamma \in \mathcal{F}_h^{ID}} \left(\int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla u \rangle)^2 dS \right)^{1/2} \left(\int_{\Gamma} \sigma [v]^2 dS \right)^{1/2} \\ &\leq \left(\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla u \rangle)^2 dS \right)^{1/2} \left(\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma [v]^2 dS \right)^{1/2}, \end{aligned} \quad (2.115)$$

and

$$\chi_3 \leq \left(\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla v \rangle)^2 dS \right)^{1/2} \left(\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma [u]^2 dS \right)^{1/2}. \quad (2.116)$$

Using the discrete Cauchy inequality, from (2.114)–(2.116) we derive the bound

$$|a_h(u, v)| \leq |u|_{H^1(\Omega, \mathcal{T}_h)} |v|_{H^1(\Omega, \mathcal{T}_h)} \quad (2.117)$$

$$\begin{aligned}
& + \left(\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla u \rangle)^2 dS \right)^{1/2} \left(\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma [v]^2 dS \right)^{1/2} \\
& + \left(\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla v \rangle)^2 dS \right)^{1/2} \left(\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma [u]^2 dS \right)^{1/2} \\
& \leq \left(|u|_{H^1(\Omega, \mathcal{T}_h)}^2 + \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla u \rangle)^2 dS + J_h^\sigma(u, u) \right)^{1/2} \\
& \times \left(|v|_{H^1(\Omega, \mathcal{T}_h)}^2 + \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla v \rangle)^2 dS + J_h^\sigma(v, v) \right)^{1/2} \\
& = \|u\|_{1, \sigma} \|v\|_{1, \sigma}.
\end{aligned}$$

□

Exercises 2.29 Prove that $\|\cdot\|_{1, \sigma}$ introduced by (2.112) defines a norm in the broken Sobolev space $H^2(\Omega, \mathcal{T}_h)$.

Corollary 2.30 By virtue of (2.47a) and (2.47b), Lemma 2.28 and Exercise 2.29, the bilinear forms A_h^s and A_h^n are bounded with respect to the norm $\|\cdot\|_{1, \sigma}$ in the broken Sobolev space $H^2(\Omega, \mathcal{T}_h)$.

Exercises 2.31 Prove Corollary 2.30.

Further, we pay attention on the expression $J_h^\sigma(u, v)$ for $u, v \in H^1(\Omega, \mathcal{T}_h)$.

Lemma 2.32 Let assumptions (2.104), (2.19) and (2.20) be satisfied. Then

$$|J_h^\sigma(u, v)| \leq J_h^\sigma(u, u)^{1/2} J_h^\sigma(v, v)^{1/2} \quad \forall u, v \in H^1(\Omega, \mathcal{T}_h), \quad (2.118)$$

and

$$\begin{aligned}
J_h^\sigma(v, v) & \leq \frac{2C_W C_M}{C_T} \sum_{K \in \mathcal{T}_h} \left(h_K^{-2} \|v\|_{L^2(K)}^2 + h_K^{-1} \|v\|_{L^2(K)} |v|_{H^1(K)} \right) \\
& \leq \frac{C_W C_M}{C_T} \sum_{K \in \mathcal{T}_h} \left(3h_K^{-2} \|v\|_{L^2(K)}^2 + |v|_{H^1(K)}^2 \right) \quad \forall v \in H^1(\Omega, \mathcal{T}_h).
\end{aligned} \quad (2.119)$$

Proof Let $u, v \in H^1(\Omega, \mathcal{T}_h)$. By the definition (2.41) of the form J_h^σ and the Cauchy inequality,

$$\begin{aligned}
|J_h^\sigma(u, v)| &\leq \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma[u][v] dS \\
&\leq \left(\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma[u]^2 dS \right)^{1/2} \left(\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma[v]^2 dS \right)^{1/2} \\
&= J_h^\sigma(u, u)^{1/2} J_h^\sigma(v, v)^{1/2}.
\end{aligned} \tag{2.120}$$

Further, the definition of the form J_h^σ , (2.104), (2.20) and (2.108) imply that

$$J_h^\sigma(v, v) = \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma[v]^2 dS = \sum_{\Gamma \in \mathcal{F}_h^{ID}} \frac{C_W}{h_{\Gamma}} \|[v]^2\|_{L^2(\Gamma)} \leq \frac{2C_W}{C_T} \sum_{K \in \mathcal{T}_h} h_K^{-1} \|v\|_{L^2(\partial K)}^2.$$

Now, using the multiplicative trace inequality (2.78), we get

$$J_h^\sigma(v, v) \leq \frac{2C_W C_M}{C_T} \sum_{K \in \mathcal{T}_h} \left(h_K^{-2} \|v\|_{L^2(K)}^2 + h_K^{-1} \|v\|_{L^2(K)} |v|_{H^1(K)} \right). \tag{2.121}$$

The last relation in (2.119) follows from (2.121) and the Young inequality. \square

Lemmas 2.28 and 2.32 immediately imply the boundedness also of the forms $A_h^{s,\sigma}$, $A_h^{n,\sigma}$ and $A_h^{i,\sigma}$ with respect to the norm $\|\cdot\|_{1,\sigma}$.

Corollary 2.33 *Let assumptions (2.104), (2.19) and (2.20) be satisfied. Then the forms A_h defined by (2.47) satisfy the estimate*

$$|A_h(u, v)| \leq 2\|u\|_{1,\sigma} \|v\|_{1,\sigma} \quad \forall u, v \in H^2(\Omega, \mathcal{T}_h). \tag{2.122}$$

Proof For the boundedness of $A_h = A_h^s$ and $A_h = A_h^n$, see Corollary (2.30). Let $A_h = A_h^{s,\sigma}$ or $A_h = A_h^{n,\sigma}$ or $A_h = A_h^{i,\sigma}$. Then, by virtue of (2.47c)–(2.47e), Lemmas 2.28 and 2.32 we have

$$\begin{aligned}
|A_h(u, v)| &\leq |a_h(u, v)| + |J_h^\sigma(u, v)| \leq \|u\|_{1,\sigma} \|v\|_{1,\sigma} + J_h^\sigma(u, u)^{1/2} J_h^\sigma(v, v)^{1/2} \\
&\leq \|u\|_{1,\sigma} \|v\|_{1,\sigma} + \|u\|_{1,\sigma} \|v\|_{1,\sigma} = 2\|u\|_{1,\sigma} \|v\|_{1,\sigma}.
\end{aligned}$$

\square

The following lemma allows us to estimate the expressions with integrals over $\Gamma \in \mathcal{F}_h$ in terms of norms over elements $K \in \mathcal{T}_h$.

Lemma 2.34 *Let the weight σ be defined by (2.104). Then, under assumptions (2.19) and (2.20), for any $v \in H^2(\Omega, \mathcal{T}_h)$ the following estimate holds:*

$$\begin{aligned}
\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla v \rangle)^2 dS &\leq \frac{C_G C_M}{C_W} \sum_{K \in \mathcal{T}_h} \left(h_K \|\nabla v\|_{L^2(K)} |\nabla v|_{H^1(K)} + \|\nabla v\|_{L^2(K)}^2 \right) \\
&= \frac{C_G C_M}{C_W} \sum_{K \in \mathcal{T}_h} \left(h_K |v|_{H^1(K)} |v|_{H^2(K)} + |v|_{H^1(K)}^2 \right) \\
&\leq \frac{C_G C_M}{2C_W} \sum_{K \in \mathcal{T}_h} \left(h_K^2 |v|_{H^2(K)}^2 + 3|v|_{H^1(K)}^2 \right). \quad (2.123)
\end{aligned}$$

Moreover, if $v \in S_{hp}$, then

$$\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla v_h \rangle)^2 dS \leq \frac{C_G C_M}{C_W} (C_I + 1) |v_h|_{H^1(\Omega, \mathcal{T}_h)}^2. \quad (2.124)$$

Proof Using (2.109) and the multiplicative trace inequality (2.78), we find that

$$\begin{aligned}
\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla v \rangle)^2 dS \\
\leq \frac{C_G}{C_W} \sum_{K \in \mathcal{T}_h} h_K \|\nabla v\|_{L^2(\partial K)}^2 \\
\leq \frac{C_G C_M}{C_W} \sum_{K \in \mathcal{T}_h} h_K \left(\|\nabla v\|_{L^2(K)} |\nabla v|_{H^1(K)} + h_K^{-1} \|\nabla v\|_{L^2(K)}^2 \right),
\end{aligned}$$

which is the first inequality in (2.123). The second one directly follows from the Young inequality.

If $v \in S_{hp}$, then (2.123) and the inverse inequality (2.86) imply that

$$\begin{aligned}
\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla v_h \rangle)^2 dS &\leq \frac{C_G C_M}{C_W} \sum_{K \in \mathcal{T}_h} \left(C_I \|\nabla v_h\|_{L^2(K)}^2 + \|\nabla v_h\|_{L^2(K)}^2 \right) \\
&= \frac{C_G C_M}{C_W} (C_I + 1) \sum_{K \in \mathcal{T}_h} \|\nabla v_h\|_{L^2(K)}^2 = \frac{C_G C_M}{C_W} (C_I + 1) |v_h|_{H^1(\Omega, \mathcal{T}_h)}^2,
\end{aligned}$$

which we wanted to prove. \square

We continue in the derivation of various inequalities based on the estimation of the $\|\cdot\|_{1,\sigma}$ -norm.

Lemma 2.35 *Under assumptions of Lemma 2.34, there exist constants C_σ , $\tilde{C}_\sigma > 0$ such that*

$$J_h^\sigma(u, u)^{1/2} \leq \|u\| \leq \|u\|_{1,\sigma} \leq C_\sigma R_a(u) \quad \forall u \in H^2(\Omega, \mathcal{T}_h), \quad h \in (0, \bar{h}), \quad (2.125)$$

$$J_h^\sigma(v_h, v_h)^{1/2} \leq |||v_h||| \leq \|v_h\|_{1,\sigma} \leq \tilde{C}_\sigma |||v_h||| \quad \forall v_h \in S_{hp}, \quad h \in (0, \bar{h}), \quad (2.126)$$

where

$$R_a(u) = \left(\sum_{K \in \mathcal{T}_h} \left(|u|_{H^1(K)}^2 + h_K^2 |u|_{H^2(K)}^2 + h_K^{-2} \|u\|_{L^2(K)}^2 \right) \right)^{1/2}, \quad u \in H^2(\Omega, \mathcal{T}_h). \quad (2.127)$$

Proof The first two inequalities in (2.125) as well as in (2.126) follow immediately from the definition of the DG-norm (2.103) and the $\|\cdot\|_{1,\sigma}$ -norm (2.112). Moreover, in view of (2.123) and (2.119), for $u \in H^2(\Omega, \mathcal{T}_h)$ we have

$$\begin{aligned} \|u\|_{1,\sigma}^2 &= |u|_{H^1(\Omega, \mathcal{T}_h)}^2 + J_h^\sigma(u, u) + \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla u \rangle)^2 dS \\ &\leq \sum_{K \in \mathcal{T}_h} |u|_{H^1(K)}^2 + \frac{C_W C_M}{C_T} \sum_{K \in \mathcal{T}_h} \left(3h_K^{-2} \|u\|_{L^2(K)}^2 + |u|_{H^1(K)}^2 \right) \\ &\quad + \frac{C_G C_M}{2C_W} \sum_{K \in \mathcal{T}_h} \left(h_K^2 |u|_{H^2(K)}^2 + 3|u|_{H^1(K)}^2 \right). \end{aligned}$$

Now, after a simple manipulation, we get

$$\begin{aligned} \|u\|_{1,\sigma}^2 &\leq \sum_{K \in \mathcal{T}_h} \left(|u|_{H^1(K)}^2 \left(1 + \frac{3C_G C_M}{2C_W} + \frac{C_W C_M}{C_T} \right) \right. \\ &\quad \left. + |u|_{H^2(K)}^2 h_K^2 \frac{C_G C_M}{2C_W} + \|u\|_{L^2(K)}^2 h_K^{-2} \frac{3C_W C_M}{C_T} \right). \end{aligned}$$

Hence, (2.125) holds with

$$C_\sigma = \left(\max \left(1 + \frac{3C_G C_M}{2C_W} + \frac{C_W C_M}{C_T}, \frac{C_G C_M}{2C_W}, \frac{3C_W C_M}{C_T} \right) \right)^{1/2}.$$

Further, if $v_h \in S_{hp}$, then (2.112), (2.124) and (2.103) immediately imply (2.126) with $\tilde{C}_\sigma = (1 + C_G C_M(C_I + 1)/C_W)^{1/2}$. \square

In what follows, we are concerned with properties of the bilinear forms A_h defined by (2.47). First, we prove the continuity of the bilinear forms A_h defined by (2.47) in the space S_{hp} with respect to the norm $||| \cdot |||$.

Lemma 2.36 *Let assumptions (2.104), (2.19) and (2.20) be satisfied. Then there exists a constant $C_B > 0$ such that the form A_h defined by (2.47) satisfies the estimate*

$$|A_h(u_h, v_h)| \leq C_B |||u_h||| |||v_h||| \quad \forall u_h, v_h \in S_{hp}. \quad (2.128)$$

Proof Estimates (2.122) and (2.126) give (2.128) with $C_B = 2\tilde{C}_\sigma^2$. \square

Further, we prove an inequality similar to (2.128) replacing $u_h \in S_{hp}$ by $u \in H^2(\Omega, \mathcal{T}_h)$.

Lemma 2.37 *Let assumptions (2.19), (2.20) and (2.104) be satisfied. Then there exists a constant $\tilde{C}_B > 0$ such that*

$$|A_h(u, v_h)| \leq \tilde{C}_B R_a(u) |||v_h||| \quad \forall u \in H^2(\Omega, \mathcal{T}_h) \quad \forall v_h \in S_{hp} \quad \forall h(0, \bar{h}), \quad (2.129)$$

where R_a is defined by (2.127).

Proof By (2.122) and (2.125),

$$|A_h(u, v_h)| \leq 2\|u\|_{1,\sigma} \|v_h\|_{1,\sigma} \leq 2C_\sigma \tilde{C}_\sigma R_a(u) |||v_h|||,$$

which is (2.129) with $\tilde{C}_B = 2C_\sigma \tilde{C}_\sigma$. \square

2.6.3 Coercivity of Diffusion Bilinear Forms

Lemma 2.38 (NIPG coercivity) *For any $C_W > 0$ the bilinear form $A_h^{n,\sigma}$ defined by (2.47d) satisfies the coercivity condition*

$$A_h^{n,\sigma}(v, v) \geq |||v|||^2 \quad \forall v \in H^2(\Omega, \mathcal{T}_h). \quad (2.130)$$

Proof From (2.45b) and (2.47d) it immediately follows that

$$A_h^{n,\sigma}(v, v) = a_h^n(v, v) + J_h^\sigma(v, v) = |v|_{H^1(\Omega, \mathcal{T}_h)}^2 + J_h^\sigma(v, v) = |||v|||^2, \quad (2.131)$$

which we wanted to prove. \square

The proof of the coercivity of the symmetric bilinear form $A_h^{s,\sigma}$ is more complicated.

Lemma 2.39 (SIPG coercivity) *Let assumptions (2.19) and (2.20) be satisfied, let*

$$C_W \geq 4C_G C_M (1 + C_I), \quad (2.132)$$

where C_M , C_I and C_G are the constants from (2.78), (2.86) and (2.20), respectively, and let the penalty parameter σ be given by (2.104) for all $\Gamma \in \mathcal{F}_h^{ID}$. Then

$$A_h^{s,\sigma}(v_h, v_h) \geq \frac{1}{2} |||v_h|||^2 \quad \forall v_h \in S_{hp} \quad \forall h \in (0, \bar{h}).$$

Proof Let $\delta > 0$. Then from (2.41), (2.104), (2.45a) and the Cauchy and Young inequalities it follows that

$$\begin{aligned}
 a_h^s(v_h, v_h) &= |v_h|_{H^1(\Omega, \mathcal{T}_h)}^2 - 2 \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \mathbf{n} \cdot \langle \nabla v_h \rangle [v_h] \, dS \\
 &\geq |v_h|_{H^1(\Omega, \mathcal{T}_h)}^2 - 2 \left\{ \frac{1}{\delta} \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} h_{\Gamma} (\mathbf{n} \cdot \langle \nabla v_h \rangle)^2 \, dS \right\}^{\frac{1}{2}} \left\{ \delta \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \frac{1}{h_{\Gamma}} [v_h]^2 \, dS \right\}^{\frac{1}{2}} \\
 &\geq |v_h|_{H^1(\Omega, \mathcal{T}_h)}^2 - \omega - \frac{\delta}{C_W} J_h^{\sigma}(v_h, v_h),
 \end{aligned} \tag{2.133}$$

where

$$\omega = \frac{1}{\delta} \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} h_{\Gamma} |\langle \nabla v_h \rangle|^2 \, dS. \tag{2.134}$$

Further, from assumption (2.20), inequality (2.107), the multiplicative trace inequality (2.78) and the inverse inequality (2.86) we get

$$\begin{aligned}
 \omega &\leq \frac{C_G}{\delta} \sum_{K \in \mathcal{T}_h} h_K \|\nabla v_h\|_{L^2(\partial K)}^2 \\
 &\leq \frac{C_G C_M}{\delta} \sum_{K \in \mathcal{T}_h} h_K \left(|v_h|_{H^1(K)} |\nabla v_h|_{H^1(K)} + h_K^{-1} |v_h|_{H^1(K)}^2 \right) \\
 &\leq \frac{C_G C_M (1 + C_I)}{\delta} |v_h|_{H^1(\Omega, \mathcal{T}_h)}^2.
 \end{aligned} \tag{2.135}$$

Now let us choose

$$\delta = 2C_G C_M (1 + C_I). \tag{2.136}$$

Then it follows from (2.132) and (2.133)–(2.136) that

$$\begin{aligned}
 a_h^s(v_h, v_h) &\geq \frac{1}{2} \left(|v_h|_{H^1(\Omega, \mathcal{T}_h)}^2 - \frac{4C_G C_M (1 + C_I)}{C_W} J_h^{\sigma}(v_h, v_h) \right) \\
 &\geq \frac{1}{2} \left(|v_h|_{H^1(\Omega, \mathcal{T}_h)}^2 - J_h^{\sigma}(v_h, v_h) \right).
 \end{aligned} \tag{2.137}$$

Finally, definition (2.47c) of the form $A_h^{s, \sigma}$ and (2.137) imply that

$$\begin{aligned}
 A_h^{s, \sigma}(v_h, v_h) &= a_h^s(v_h, v_h) + J_h^{\sigma}(v_h, v_h) \\
 &\geq \frac{1}{2} \left(|v_h|_{H^1(\Omega, \mathcal{T}_h)}^2 + J_h^{\sigma}(v_h, v_h) \right) = \frac{1}{2} \|v_h\|^2,
 \end{aligned} \tag{2.138}$$

which we wanted to prove. \square

Lemma 2.40 (IIPG coercivity) *Let assumptions (2.19) and (2.20) be satisfied, let*

$$C_W \geq C_G C_M (1 + C_I), \quad (2.139)$$

where C_M , C_I and C_G are constants from (2.78), (2.86) and (2.20), respectively, and let the penalty parameter σ be given by (2.104) for all $\Gamma \in \mathcal{F}_h^{ID}$. Then

$$A_h^{i,\sigma}(v_h, v_h) \geq \frac{1}{2} \|v_h\|^2 \quad \forall v_h \in S_{hp}.$$

Proof The proof is almost identical with the proof of the previous lemma. □

Corollary 2.41 *We can summarize the above results in the following way. We have*

$$A_h(v_h, v_h) \geq C_C \|v_h\|^2 \quad \forall v_h \in S_{hp}, \quad (2.140)$$

with

$$\begin{aligned} C_C &= 1 && \text{for } A_h = A_h^{n,\sigma} && \text{if } C_W > 0, \\ C_C &= 1/2 && \text{for } A_h = A_h^{s,\sigma} && \text{if } C_W \geq 4C_G C_M (1 + C_I), \\ C_C &= 1/2 && \text{for } A_h = A_h^{i,\sigma} && \text{if } C_W \geq C_G C_M (1 + C_I). \end{aligned}$$

Corollary 2.42 *By virtue of Corollary 1.7, the coercivity of the forms A_h implies the existence and uniqueness of the solution of the discrete problems (2.49c)–(2.49e) (SIPG, NIPG and IIPG method).*

2.7 Error Estimates

In this section, we derive error estimates of the SIPG, NIPG and IIPG variants of the DGM applied to the numerical solution of the Poisson problem (2.1). Namely, the error $u_h - u$ will be estimated in the DG-norm and the $L^2(\Omega)$ -norm.

2.7.1 Estimates in the DG-Norm

Let $u \in H^2(\Omega)$ denote the exact strong solution of problem (2.1) and let $u_h \in S_{hp}$ be the approximate solution obtained by method (2.54), where the forms A_h and ℓ_h are defined by (2.47c)–(2.47e) and (2.48c)–(2.48e), respectively. The error of the method is defined as the function $e_h = u_h - u \in H^2(\Omega, \mathcal{T}_h)$. It can be written in the form

$$e_h = \xi + \eta, \quad \text{with } \xi = u_h - \Pi_{hp} u \in S_{hp}, \quad \eta = \Pi_{hp} u - u \in H^2(\Omega, \mathcal{T}_h), \quad (2.141)$$

where Π_{hp} is the S_{hp} -interpolation defined by (2.90). Hence, we split the error into two parts ξ and η . The term η represents the error of the S_{hp} -interpolation of the function u . (It is possible to say that η approximates the *distance* of the exact solution from the space S_{hp} , where the approximate solution is sought.) The term η can be simply estimated on the basis of the approximation properties (2.92) and (2.97). On the other hand, the term ξ represents the *distance* between the approximate solution u_h and the projection of the exact solution on the space S_{hp} . The estimation of ξ is sometimes more complicated.

We suppose that the system of triangulations $\{\mathcal{T}_h\}_{h \in (0, \bar{h})}$ satisfies the shape-regularity assumption (2.19) and that the equivalence condition (2.20) holds.

First, we prove the so-called *abstract error estimate*, representing a bound of the error in terms of the S_{hp} -interpolation error η .

Theorem 2.43 *Let assumptions (2.19) and (2.20) be satisfied and let the exact solution of problem (2.1) satisfy the condition $u \in H^2(\Omega)$. Then there exists a constant $C_{AE} > 0$ such that*

$$\|e_h\| \leq C_{AE} R_a(\eta) = C_{AE} R_a(\Pi_{hp}u - u), \quad h \in (0, \bar{h}), \quad (2.142)$$

where $R_a(\eta)$ is given by (2.127).

Proof We express the error by (2.141), i.e., $e_h = u_h - u = \xi + \eta$. The error e_h satisfies the Galerkin orthogonality condition (2.57), which is equivalent to the relation

$$A_h(\xi, v_h) = -A_h(\eta, v_h) \quad \forall v_h \in S_{hp}. \quad (2.143)$$

If we set $v_h := \xi \in S_{hp}$ in (2.143) and use (2.47c)–(2.47e) and the coercivity (2.140), we find that

$$C_C \|\xi\|^2 \leq A_h(\xi, \xi) = -A_h(\eta, \xi). \quad (2.144)$$

Now we apply Lemma 2.37 and get

$$|A_h(\eta, \xi)| \leq \tilde{C}_B R_a(\eta) \|\xi\|.$$

The above and (2.144) already imply that

$$\|\xi\| \leq \frac{\tilde{C}_B}{C_C} R_a(\eta). \quad (2.145)$$

Obviously,

$$\|e_h\| \leq \|\xi\| + \|\eta\|. \quad (2.146)$$

Finally, (2.125) gives

$$|||\eta||| \leq C_\sigma R_a(\eta). \quad (2.147)$$

Hence, (2.146), (2.145) and (2.147) yield the abstract error estimate (2.142) with $C_{AE} = C_\sigma + \tilde{C}_B/C_C$. \square

The abstract error estimate is the basis for estimating the error e_h in terms of the mesh-size h .

Theorem 2.44 (DG-norm error estimate) *Let us assume that $s \geq 2$, $p \geq 1$, are integers, $u \in H^s(\Omega)$ is the solution of problem (2.1), $\{\mathcal{T}_h\}_{h \in (0, \bar{h})}$ is a system of triangulations of the domain Ω satisfying the shape-regularity condition (2.19), and the equivalence condition (2.20) (cf. Lemma 2.5). Moreover, let the penalty constant C_W satisfy the conditions from Corollary 2.41. Let $u_h \in S_{hp}$ be the approximate solution obtained by using of the SIPG, NIPG or IIPG method (2.49c)–(2.49e). Then the error $e_h = u_h - u$ satisfies the estimate*

$$|||e_h||| \leq C_1 h^{\mu-1} |u|_{H^\mu(\Omega)}, \quad h \in (0, \bar{h}), \quad (2.148)$$

where $\mu = \min(p+1, s)$ and C_1 is a constant independent of h and u . Hence, if $s \geq p+1$, we get the error estimate

$$|||e_h||| \leq C_1 h^p |u|_{H^{p+1}(\Omega)}.$$

Proof It is enough to use the abstract error estimate (2.142), where the expressions $|\eta|_{H^1(K)}$, $|\eta|_{H^2(K)}$ and $\|\eta\|_{L^2(K)}$, $K \in \mathcal{T}_h$, are estimated on the basis of the approximation properties (2.93)–(2.95), rewritten for $\eta|_K = (\Pi_{hp}u - u)|_K = \pi_{K,p}(u|_K) - u|_K$ and $K \in \mathcal{T}_h$:

$$\begin{aligned} \|\eta\|_{L^2(K)} &\leq C_A h_K^\mu |u|_{H^\mu(K)}, \\ |\eta|_{H^1(K)} &\leq C_A h_K^{\mu-1} |u|_{H^\mu(K)}, \\ |\eta|_{H^2(K)} &\leq C_A h_K^{\mu-2} |u|_{H^\mu(K)}. \end{aligned} \quad (2.149)$$

Thus, the inequality $h_K \leq h$ and the relation $\sum_{K \in \mathcal{T}_h} |u|_{H^\mu(K)}^2 = |u|_{H^\mu(\Omega)}^2$ imply

$$\begin{aligned} R_a(\eta) &= \left(\sum_{K \in \mathcal{T}_h} \left(|\eta|_{H^1(K)}^2 + h_K^2 |\eta|_{H^2(K)}^2 + h_K^{-2} \|\eta\|_{L^2(K)}^2 \right) \right)^{1/2} \\ &\leq \sqrt{3} C_A h^{\mu-1} |u|_{H^\mu(\Omega)}, \end{aligned} \quad (2.150)$$

which together with (2.142) gives (2.148) with the constant $C_1 = \sqrt{3} C_{AE} C_A$. \square

In order to derive an error estimate in the $L^2(\Omega)$ -norm we present the following result.

Lemma 2.45 (Broken Poincaré inequality) *Let the system $\{\mathcal{T}_h\}_{h \in (0, \bar{h})}$ of triangulations satisfy the shape-regularity assumption (2.19). Then there exists a constant $C > 0$ independent of h and v_h such that*

$$\|v_h\|_{L^2(\Omega)}^2 \leq C \left(\sum_{K \in \mathcal{T}_h} |v_h|_{H^1(K)}^2 + \sum_{\Gamma \in \mathcal{F}_h^{ID}} \frac{1}{\text{diam}(\Gamma)} \|[v_h]\|_{L^2(\Gamma)}^2 \right) \quad (2.151)$$

$$\forall v_h \in S_{hp} \quad \forall h \in (0, \bar{h}).$$

The proof of the broken Poincaré inequality (2.151) was carried out in [7] in the case where Ω is a convex polygonal domain, $\partial\Omega_D = \partial\Omega$ and the assumption (MA2) in Sect. 2.3.2 is satisfied. The proof of inequality (2.151) in a general case with the nonempty Neumann part of the boundary can be found in [36].

From Theorem 2.44 and (2.151) we obtain the following result.

Corollary 2.46 ($L^2(\Omega)$ -(suboptimal) error estimate) *Let the assumptions of Theorem 2.44 be satisfied. Then*

$$\|e_h\|_{L^2(\Omega)} \leq C_2 h^{\mu-1} |u|_{H^\mu(\Omega)}, \quad h \in (0, \bar{h}), \quad (2.152)$$

where C_2 is a constant independent of h . Hence, if $s \geq p + 1$, we get the error estimate

$$\|e_h\|_{L^2(\Omega)} \leq C_2 h^p |u|_{H^{p+1}(\Omega)}. \quad (2.153)$$

Remark 2.47 The error estimate (2.153), which is of order $O(h^p)$, is suboptimal with respect to the approximation property (2.97) with $q = 0$, $\mu = p + 1 \leq s$ of the space S_{hp} giving the order $O(h^{p+1})$. In the next section we prove an optimal error estimate in the $L^2(\Omega)$ -norm for SIPG method using the Aubin–Nitsche technique.

2.7.2 Optimal $L^2(\Omega)$ -Error Estimate

Our further aim is to derive the optimal error estimate in the $L^2(\Omega)$ -norm. It will be based on the *duality technique* sometimes called the *Aubin–Nitsche trick*. Since this approach requires the symmetry of the corresponding bilinear form and the regularity of the exact solution to the dual problem, we consider the SIPG method applied to problem (2.1) with $\partial\Omega_D = \partial\Omega$ and $\partial\Omega_N = \emptyset$. This means that we seek u satisfying

$$-\Delta u = f \quad \text{in } \Omega, \quad (2.154a)$$

$$u = u_D \quad \text{on } \partial\Omega. \quad (2.154b)$$

Moreover, for an arbitrary $z \in L^2(\Omega)$, we consider the *dual problem*: Given $z \in L^2(\Omega)$, find ψ such that

$$-\Delta\psi = z \quad \text{in } \Omega, \quad \psi = 0 \quad \text{on } \partial\Omega. \quad (2.155)$$

Under the notation

$$V = H_0^1(\Omega) = \left\{ v \in H^1(\Omega); v = 0 \text{ on } \partial\Omega \right\}, \quad (2.156)$$

the weak formulation of (2.155) reads: Find $\psi \in V$ such that

$$\int_{\Omega} \nabla\psi \cdot \nabla v \, dx = \int_{\Omega} zv \, dx = (z, v)_{L^2(\Omega)} \quad \forall v \in V. \quad (2.157)$$

Let us assume that $\psi \in H^2(\Omega)$ and that there exists a constant $C_D > 0$, independent of z , such that

$$\|\psi\|_{H^2(\Omega)} \leq C_D \|z\|_{L^2(\Omega)}. \quad (2.158)$$

This is true provided the polygonal (polyhedral) domain Ω is convex, as follows from [153]. (See Remark 2.50.) Let us note that $H^2(\Omega) \subset C(\overline{\Omega})$, if $d \leq 3$.

Let A_h be the symmetric bilinear form given by (2.47c), i.e.,

$$A_h(u, v) = a_h^s(u, v) + J_h^\sigma(u, v), \quad u, v \in H^2(\Omega, \mathcal{T}_h), \quad (2.159)$$

where a_h^s and J_h^σ are defined by (2.45a) and (2.105), respectively.

First, we prove the following auxiliary result.

Lemma 2.48 *Let $\psi \in H^2(\Omega)$ be the solution of problem (2.155). Then*

$$A_h(\psi, v) = (v, z)_{L^2(\Omega)} \quad \forall v \in H^2(\Omega, \mathcal{T}_h). \quad (2.160)$$

Proof The function $\psi \in H^2(\Omega)$ satisfies the conditions

$$[\psi]_\Gamma = 0 \quad \forall \Gamma \in \mathcal{F}_h^I, \quad \psi|_{\partial\Omega} = 0. \quad (2.161)$$

Let $v \in H^2(\Omega, \mathcal{T}_h)$. Using (2.155), (2.161) and Green's theorem, we obtain

$$\begin{aligned} (v, z)_{L^2(\Omega)} &= \int_{\Omega} zv \, dx = - \int_{\Omega} \Delta\psi v \, dx \\ &= \sum_{K \in \mathcal{T}_h} \int_K \nabla\psi \cdot \nabla v \, dx - \sum_{K \in \mathcal{T}_h} \int_{\partial K} \nabla\psi \cdot \mathbf{n} v \, dS \\ &= \sum_{K \in \mathcal{T}_h} \int_K \nabla\psi \cdot \nabla v \, dx \end{aligned}$$

$$\begin{aligned}
& - \left(\sum_{\Gamma \in \mathcal{F}_h^I} \int_{\Gamma} \langle \nabla \psi \rangle \cdot \mathbf{n} [v] \, dS + \sum_{\Gamma \in \mathcal{F}_h^I} \int_{\Gamma} \langle \nabla v \rangle \cdot \mathbf{n} [\psi] \, dS \right) \\
& - \left(\sum_{\Gamma \in \mathcal{F}_h^B} \int_{\Gamma} \nabla \psi \cdot \mathbf{n} v \, dS + \sum_{\Gamma \in \mathcal{F}_h^B} \int_{\Gamma} \nabla v \cdot \mathbf{n} \psi \, dS \right) \\
& + \left(\sum_{\Gamma \in \mathcal{F}_h^I} \int_{\Gamma} \sigma [\psi] [v] \, dS + \sum_{\Gamma \in \mathcal{F}_h^B} \int_{\Gamma} \sigma \psi v \, dS \right).
\end{aligned}$$

Hence, in view of the definition of the form A_h , we have (2.160). \square

Theorem 2.49 ($L^2(\Omega)$ -optimal error estimate) *Let us assume that $s \geq 2$, $p \geq 1$, are integers, Ω is a bounded convex polyhedral domain, $u \in H^s(\Omega)$ is the solution of problem (2.1), $\{\mathcal{T}_h\}_{h \in (0, \bar{h})}$ is a system of triangulations of the domain Ω satisfying the shape-regularity condition (2.19), and the equivalence condition (2.20) (cf. Lemma 2.5). Moreover, let the penalty constant C_W satisfy the condition from Corollary 2.41. Let $u_h \in S_{hp}$ be the approximate solution obtained using the SIPG method (2.49c) (i.e., $\Theta = 1$ and the form $A_h = A_h^{\sigma, s}$ is given by (2.45a) and (2.47c). Then*

$$\|e_h\|_{L^2(\Omega)} \leq C_3 h^\mu |u|_{H^\mu(\Omega)}, \quad (2.162)$$

where $e_h = u_h - u$, $\mu = \min\{p + 1, s\}$ and C_3 is a constant independent of h and u .

Proof Let $\psi \in H^2(\Omega)$ be the solution of the dual problem (2.157) with $z := e_h = u_h - u \in L^2(\Omega)$ and let $\Pi_{h1}\psi \in S_{h1}$ be the approximation of ψ defined by (2.90) with $p = 1$. By (2.160), we have

$$A_h(\psi, v) = (e_h, v)_{L^2(\Omega)} \quad \forall v \in H^2(\Omega, \mathcal{T}_h). \quad (2.163)$$

The symmetry of the form A_h , the Galerkin orthogonality (2.57) of the error and (2.163) with $v := e_h$ yield

$$\begin{aligned}
\|e_h\|_{L^2(\Omega)}^2 &= A_h(\psi, e_h) = A_h(e_h, \psi) \\
&= A_h(e_h, \psi - \Pi_{h1}\psi).
\end{aligned} \quad (2.164)$$

Moreover, from (2.122), it follows that

$$A_h(e_h, \psi - \Pi_{h1}\psi) \leq 2\|e_h\|_{1,\sigma} \|\psi - \Pi_{h1}\psi\|_{1,\sigma}, \quad (2.165)$$

where, by (2.112),

$$\|v\|_{1,\sigma}^2 = \|v\|^2 + \sum_{\Gamma \in \mathcal{T}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla v \rangle)^2 dS. \quad (2.166)$$

By (2.125) and (2.150) (with $\mu = 2$), we have

$$\|\psi - \Pi_{h1}\psi\|_{1,\sigma} \leq C_{\sigma} R_a(\psi - \Pi_{h1}\psi) \leq \sqrt{3} C_{\sigma} C_A h |\psi|_{H^2(\Omega)}. \quad (2.167)$$

Now, the inverse inequality (2.86) and estimates (2.100), (2.99) imply that

$$\begin{aligned} |\nabla e_h|_{H^1(K)} &= |\nabla(u - u_h)|_{H^1(K)} \\ &\leq |\nabla(u - \Pi_{hp}u)|_{H^1(K)} + |\nabla(\Pi_{hp}u - u_h)|_{H^1(K)} \\ &\leq |u - \Pi_{hp}u|_{H^2(K)} + C_I h_K^{-1} \|\nabla(\Pi_{hp}u - u_h)\|_{L^2(K)} \\ &\leq C_A h_K^{\mu-2} |u|_{H^{\mu}(K)} + C_I h_K^{-1} \left(\|\nabla(\Pi_{hp}u - u)\|_{L^2(K)} + \|\nabla(u - u_h)\|_{L^2(K)} \right) \\ &\leq C_A (1 + C_I) h_K^{\mu-2} |u|_{H^{\mu}(K)} + C_I h_K^{-1} \|\nabla e_h\|_{L^2(K)}. \end{aligned} \quad (2.168)$$

By (2.123), (2.168) and the discrete Cauchy inequality,

$$\begin{aligned} \sum_{\Gamma \in \mathcal{T}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla e_h \rangle)^2 dS \\ \leq \frac{C_G C_M}{C_W} \sum_{K \in \mathcal{T}_h} \left(h_K \|\nabla e_h\|_{L^2(K)} |\nabla e_h|_{H^1(K)} + \|\nabla e_h\|_{L^2(K)}^2 \right) \\ \leq \frac{C_G C_M}{C_W} \left\{ C_A (1 + C_I) h^{\mu-1} |e_h|_{H^1(\Omega, \mathcal{T}_h)} |u|_{H^{\mu}(\Omega)} + (1 + C_I) |e_h|_{H^1(\Omega, \mathcal{T}_h)}^2 \right\}. \end{aligned} \quad (2.169)$$

Since $|e_h|_{H^1(\Omega, \mathcal{T}_h)} \leq \|e_h\|$, using (2.148) and (2.169), we have

$$\sum_{\Gamma \in \mathcal{T}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \langle \nabla e_h \rangle)^2 dS \leq \frac{C_G C_M}{C_W} C_1 (1 + C_I) (C_1 + C_A) h^{2(\mu-1)} |u|_{H^{\mu}(\Omega)}^2.$$

Thus, (2.148) and (2.166) yield the estimate

$$\|e_h\|_{1,\sigma}^2 \leq C_5 h^{2(\mu-1)} |u|_{H^{\mu}(\Omega)}^2 \quad (2.170)$$

with $C_5 = C_1 \left\{ 1 + C_G C_M C_W^{-1} (1 + C_I) (C_1 + C_A) \right\}$. It follows from (2.165), (2.167), and (2.170) that

$$A_h(e_h, \psi - \Pi_{h1}\psi) \leq C_6 h^{\mu} |\psi|_{H^2(\Omega)} |u|_{H^{\mu}(\Omega)}, \quad (2.171)$$

where $C_6 = 2\sqrt{3}C_\sigma C_A \sqrt{C_5}$.

Finally, by (2.164), (2.171), and (2.158) with $z = e_h$,

$$\|e_h\|_{L^2(\Omega)}^2 \leq C_D C_6 h^\mu \|u\|_{H^\mu(\Omega)} \|e_h\|_{L^2(\Omega)}, \quad (2.172)$$

which already implies estimate (2.162) with $C_3 = C_D C_6$. \square

Remark 2.50 As we see from the above results, if the exact solution $u \in H^{p+1}(\Omega)$ and the finite elements of degree p are used, the error is of the optimal order $O(h^{p+1})$ in the $L^2(\Omega)$ -norm. In the case, when the polygonal domain is not convex and/or the Neumann and Dirichlet parts of the boundary $\Omega_N \neq \emptyset$ and $\Omega_D \neq \emptyset$, the exact solution ψ of the dual problem (2.155) is not an element of the space $H^2(\Omega)$. Then it is necessary to work in the Sobolev–Slobodetskii spaces of functions with *noninteger derivatives* and the error in the $L^2(\Omega)$ -norm is not of the optimal order $O(h^{p+1})$. The analysis of error estimates for the DG discretization of boundary value problems with boundary singularities is the subject of works [137, 284], where optimal error estimates were obtained with the aid of a suitable graded mesh refinement. The main tools are here the Sobolev–Slobodetskii spaces and weighted Sobolev spaces. For the definitions and properties of these spaces, see [37, 209].

Remark 2.51 In [240] the Neumann problem (i.e., $\partial\Omega = \partial\Omega_N$) was solved by the NIPG approach, where the penalty coefficient σ was chosen in the form

$$\sigma|_\Gamma = \frac{C_W}{h_\Gamma^\beta}, \quad \Gamma \in \mathcal{T}_h, \quad (2.173)$$

instead of (2.104), where $\beta \geq 1/2$. If triangular grids do not contain any hanging nodes (i.e., the triangulations \mathcal{T}_h are conforming), then an optimal error estimate in the $L^2(\Omega)$ -norm of this analogue of the NIPG method was proven provided that $\beta \geq 3$ for $d = 2$ and $\beta \geq 3/2$ for $d = 3$. In this case the interior penalty is so strong that the DG methods behave like the standard conforming (i.e., continuous) finite element schemes. On the other hand, the stronger penalty causes worse computational properties of the corresponding algebraic system, see [41].

2.8 Baumann–Oden Method

In this section we analyze the Baumann–Oden scheme (2.49b). Hence, we seek $u_h \in S_{hp}$ such that

$$A_h(u_h, v_h) = \ell_h(v_h) \quad \forall v_h \in S_{hp}, \quad (2.174)$$

where $A_h(\cdot, \cdot)$ and ℓ_h are given by (2.47b) and (2.48b), respectively:

$$A_h(u, v) = \sum_{K \in \mathcal{T}_h} \int_K \nabla u \cdot \nabla v \, dx - \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} (\mathbf{n} \cdot \langle \nabla u \rangle [v] - \mathbf{n} \cdot \langle \nabla v \rangle [u]) \, dS, \quad (2.175)$$

$$\ell_h(v) = \int_{\Omega} f v \, dx + \sum_{\Gamma \in \mathcal{F}_h^N} \int_{\Gamma} g_N v \, dS + \sum_{\Gamma \in \mathcal{F}_h^D} \int_{\Gamma} (\mathbf{n} \cdot \nabla v) u_D \, dS.$$

Obviously, (2.175) gives

$$A_h(v, v) \geq |v|_{H^1(\Omega, \mathcal{T}_h)}^2 \quad \forall v \in H^2(\Omega, \mathcal{T}_h), \quad (2.176)$$

where only a seminorm stands on the right-hand side. We speak about a *weak coercivity*. (The above inequality is valid with the sign = of course.) Therefore, it is possible to derive error estimates in a seminorm only.

This method was presented and analyzed for one-dimensional diffusion problem in [12]. In [239], Rivière, Wheeler, and Girault showed how to obtain error estimates under the assumption that the polynomial degree $p \geq 2$ and the mesh is conforming. The analysis carried out in [239, Lemma 5.1] is based on the existence of an interpolation operator $I_{hp} : H^2(\Omega, \mathcal{T}_h) \rightarrow S_{hp}$ for $p \geq 2$ such that

$$\int_{\Gamma} \langle \nabla(v - I_{hp}v) \rangle \cdot \mathbf{n} \, dS = 0 \quad \forall \Gamma \in \mathcal{F}_h, \quad v \in H^2(\Omega, \mathcal{T}_h), \quad (2.177)$$

$$|I_{hp}v - v|_{H^q(\Omega, \mathcal{T}_h)} \leq \bar{C}_A h^{\mu-q} |v|_{H^{\mu}(\Omega, \mathcal{T}_h)}, \quad v \in H^s(\Omega, \mathcal{T}_h), \quad h \in (0, \bar{h}), \quad (2.178)$$

where $\mu = \min(p+1, s)$, $s \geq 2$, $q = 0, 1, 2$ and \bar{C}_A is a constant.

In the following, we present the error estimate for the Baumann–Oden method. The proof differs from the technique in [239].

Theorem 2.52 *Let $u \in H^s(\Omega)$ with $s \geq 2$ be the exact solution of problem (2.1). Let the system of triangulations $\{\mathcal{T}_h\}_{h \in (0, \bar{h})}$ satisfy the shape-regularity assumption (2.19) and the conformity assumption (MA4) from Sect. 2.3.2, and let $u_h \in S_{hp}$, $p \geq 2$, be the approximate solution given by (2.174). Then there exists a constant $C_{BO} > 0$ independent of $h \in (0, \bar{h})$ and u , such that*

$$|u - u_h|_{H^1(\Omega, \mathcal{T}_h)} \leq C_{BO} h^{\mu-1} |u|_{H^{\mu}(\Omega)}. \quad (2.179)$$

Proof Let I_{hp} be the interpolation operator satisfying (2.177) and (2.178). We put $\eta = I_{hp}u - u \in H^1(\Omega, \mathcal{T}_h)$ and $\xi = u_h - I_{hp}u \in S_{hp}$. Then $e_h = u_h - u = \eta + \xi$. From the definition (2.175) of the form A_h and the Galerkin orthogonality (2.57), we have

$$\begin{aligned}
|\xi|_{H^1(\Omega, \mathcal{T}_h)}^2 &= |I_{hp}u - u_h|_{H^1(\Omega, \mathcal{T}_h)}^2 = A_h(I_{hp}u - u_h, I_{hp}u - u_h) \quad (2.180) \\
&= A_h(I_{hp}u - u, I_{hp}u - u_h) = A_h(\eta, \xi).
\end{aligned}$$

Moreover, in view of (2.175) and (2.177),

$$A_h(I_{hp}u - u, v_h) = A_h(\eta, v_h) = 0 \quad \forall v_h \in S_{h0}, \quad (2.181)$$

where S_{h0} denotes the space of piecewise constant functions on \mathcal{T}_h . Hence, if Π_{h0} is the orthogonal projection of $L^2(\Omega)$ onto S_{h0} , then (2.47b), (2.111) and (2.181) imply that

$$\begin{aligned}
|A_h(\eta, \xi)| &\leq |A_h(\eta, \xi - \Pi_{h0}\xi)| + |A_h(\eta, \Pi_{h0}\xi)| \\
&\leq \|\eta\|_{1,\sigma} \|\xi - \Pi_{h0}\xi\|_{1,\sigma}, \quad (2.182)
\end{aligned}$$

where, by (2.112),

$$\|v\|_{1,\sigma}^2 = \|v\|^2 + \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1} (\mathbf{n} \cdot \nabla v)^2 dS. \quad (2.183)$$

Since $\Pi_{h0}|_K$ is constant on each $K \in \mathcal{T}_h$, obviously

$$|\xi - \Pi_{h0}\xi|_{H^1(K)} = |\xi|_{H^1(K)}, \quad K \in \mathcal{T}_h. \quad (2.184)$$

Moreover, it follows from the approximation properties (2.90) and (2.93) (with $\mu = 1$, $p = 0$) that

$$\|\xi - \Pi_{h0}\xi\|_{L^2(K)} \leq C_A h_K |\xi|_{H^1(K)}, \quad K \in \mathcal{T}_h. \quad (2.185)$$

Let $\psi \in H^1(\Omega, \mathcal{T}_h)$. Then, using the definition (2.105) of the form J_h^σ , the definition (2.104) of the weight σ , inequality (2.108), and the multiplicative trace inequality (2.78), we find that

$$\begin{aligned}
|||\psi|||^2 &= |\psi|_{H^1(\Omega, \mathcal{T}_h)}^2 + J_h^\sigma(\psi, \psi) \quad (2.186) \\
&\leq |\psi|_{H^1(\Omega, \mathcal{T}_h)}^2 + \frac{2C_W}{C_T} \sum_{K \in \mathcal{T}_h} h_K^{-1} \|\psi\|_{L^2(\partial K)}^2 \\
&\leq |\psi|_{H^1(\Omega, \mathcal{T}_h)}^2 + \frac{2C_W C_M}{C_T} \sum_{K \in \mathcal{T}_h} \left(h_K^{-2} \|\psi\|_{L^2(K)}^2 + h_K^{-1} \|\psi\|_{L^2(K)} |\psi|_{H^1(K)} \right).
\end{aligned}$$

Let us set $\psi = \xi - \Pi_{h0}\xi$ in (2.186). Then, in view of (2.184) and (2.185), we get

$$\begin{aligned}
|||\xi - \Pi_{h0}\xi|||^2 &\leq (1 + 2(1 + C_A)C_AC_W C_M/C_T) \sum_{K \in \mathcal{T}_h} |\xi|_{H^1(K)}^2 \\
&= (1 + 2(1 + C_A)C_AC_W C_M/C_T) |\xi|_{H^1(\Omega, \mathcal{T}_h)}^2.
\end{aligned} \tag{2.187}$$

Moreover, using the relation $\nabla \Pi_{h0}\xi = 0$ and (2.124), we have

$$\begin{aligned}
&\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1}(\mathbf{n} \cdot \langle \nabla(\xi - \Pi_{h0}\xi) \rangle)^2 dS \\
&= \sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1}(\mathbf{n} \cdot \langle \nabla \xi \rangle)^2 dS \leq (C_G C_M (C_I + 1)/C_W) |\xi|_{H^1(\Omega, \mathcal{T}_h)}^2.
\end{aligned} \tag{2.188}$$

Therefore, (2.183), (2.187) and (2.188) imply that

$$\|\xi - \Pi_{h0}\xi\|_{1,\sigma} \leq C_7 |\xi|_{H^1(\Omega, \mathcal{T}_h)}, \tag{2.189}$$

where $C_7 = (1 + 2(1 + C_A)C_AC_W C_M/C_T + C_G C_M (C_I + 1)/C_W)^{1/2}$.

On the other hand, if we set $\psi := \eta$ in (2.186), then by (2.178) we obtain

$$\begin{aligned}
|||\eta|||^2 &\leq \bar{C}_A^2 h^{2(\mu-1)} |u|_{H^\mu(\Omega)}^2 + 4\bar{C}_A^2 C_W C_M/C_T h^{2(\mu-1)} |u|_{H^\mu(\Omega)}^2 \\
&= \bar{C}_A^2 (1 + 4C_W C_M/C_T) h^{2(\mu-1)} |u|_{H^\mu(\Omega)}^2.
\end{aligned} \tag{2.190}$$

Similarly, inequalities (2.123) and (2.178) give

$$\begin{aligned}
&\sum_{\Gamma \in \mathcal{F}_h^{ID}} \int_{\Gamma} \sigma^{-1}(\mathbf{n} \cdot \langle \nabla \eta \rangle)^2 dS \\
&\leq \frac{C_G C_M}{C_W} \sum_{K \in \mathcal{T}_h} \left(\|\nabla \eta\|_{L^2(K)}^2 + h_K \|\nabla \eta\|_{L^2(K)} |\nabla \eta|_{H^1(K)} \right) \\
&\leq \frac{2C_G C_M \bar{C}_A^2}{C_W} h^{2(\mu-1)} |u|_{H^\mu(\Omega)}^2.
\end{aligned} \tag{2.191}$$

Then (2.183), (2.190) and (2.191) yield

$$\|\eta\|_{1,\sigma} \leq C_8 h^{\mu-1} |u|_{H^\mu(\Omega)}, \tag{2.192}$$

where $C_8 = \bar{C}_A((1 + 4C_W C_M/C_T) + 2C_G C_M/C_W)^{1/2}$.

Further, from (2.180), (2.182), (2.189) and (2.192), we have

$$|\xi|_{H^1(\Omega, \mathcal{T}_h)}^2 = |A_h(\eta, \xi)| \leq C_7 C_8 h^{\mu-1} |\xi|_{H^1(\Omega, \mathcal{T}_h)} |u|_{H^\mu(\Omega)} \tag{2.193}$$

and thus,

$$|\xi|_{H^1(\Omega, \mathcal{T}_h)} \leq C_7 C_8 h^{\mu-1} |u|_{H^\mu(\Omega)}. \quad (2.194)$$

Finally, the triangle inequality, the definition of η and ξ , (2.178), and (2.194) imply that

$$\begin{aligned} |u - u_h|_{H^1(\Omega, \mathcal{T}_h)} &= |u - I_{hp}u|_{H^1(\Omega, \mathcal{T}_h)} + |I_{hp}u - u_h|_{H^1(\Omega, \mathcal{T}_h)} \\ &\leq (\bar{C}_A + C_7 C_8) h^{\mu-1} |u|_{H^\mu(\Omega)}, \end{aligned} \quad (2.195)$$

which proves the theorem with $C_{BO} := \bar{C}_A + C_7 C_8$. \square

2.9 Numerical Examples

In this section, we demonstrate by numerical experiments the error estimates (2.148), (2.152) and (2.162). In the first example, we assume that the exact solution is sufficiently regular. We show that the use of a higher degree of polynomial approximation increases the rate of convergence of the method. In the second example, the exact solution has a singularity. Then the order of convergence does not increase with the increasing degree of the polynomial approximation used. The computational results are in agreement with theory and show that the accuracy of the method is determined by the degree of the polynomial approximation as well as the regularity of the solution.

2.9.1 Regular Solution

Let us consider the problem of finding a function $u : \Omega = (0, 1) \times (0, 1) \rightarrow \mathbb{R}$ such that

$$\begin{aligned} -\Delta u &= 8\pi^2 \sin(2\pi x_1) \sin(2\pi x_2) \quad \text{in } \Omega, \\ u &= 0 \quad \text{on } \partial\Omega. \end{aligned} \quad (2.196)$$

It is easy to verify that the exact solution of (2.196) has the form

$$u = \sin(2\pi x_1) \sin(2\pi x_2), \quad (x_1, x_2) \in \Omega. \quad (2.197)$$

Obviously, $u \in C^\infty(\bar{\Omega})$.

We investigate the *experimental order of convergence* (EOC) of the SIPG, NIPG and IIPG methods defined by (2.49c)–(2.49e). We assume that a (semi)norm $\|e_h\|$ of the computational error behaves according to the formula

$$\|e_h\| = Ch^{\text{EOC}}, \quad (2.198)$$

where $C > 0$ is a constant, $h = \max_{K \in \mathcal{T}_h} h_K$, and $\text{EOC} \in \mathbb{R}$ is the experimental order of convergence. Since the exact solution is known and therefore $\|e_h\|$ can be exactly evaluated, it is possible to evaluate EOC in the following way. Let $\|e_{h_1}\|$ and $\|e_{h_2}\|$ be computational errors of the numerical solutions obtained on two different meshes \mathcal{T}_{h_1} and \mathcal{T}_{h_2} , respectively. Then from (2.198), eliminating the constant C , we obtain

$$\text{EOC} = \frac{\log(\|e_{h_1}\|/\|e_{h_2}\|)}{\log(h_1/h_2)}. \quad (2.199)$$

Moreover, we evaluate the *global experimental order of convergence* (GEOC) from the approximation of (2.198) with the aid of the least squares method, where all computed pairs $[h, e_h]$ are taken into account simultaneously.

We used a set of four uniform triangular grids having 128, 512, 2048, and 8192 elements, shown in Fig. 2.4. The meshes consist of right-angled triangles with the diameter $h = \sqrt{2}/\sqrt{\#\mathcal{T}_h/2}$, where $\#\mathcal{T}_h$ is the number of elements of \mathcal{T}_h . EOC is evaluated according to (2.199) for all pairs of “neighbouring” grids. Tables 2.1 and 2.2 show the computational errors in the $L^2(\Omega)$ -norm and the $H^1(\Omega, \mathcal{T}_h)$ -seminorm and EOC obtained by the SIPG, NIPG and IIPG methods using the P_p , $p = 1, \dots, 6$, polynomial approximations. These results are also visualized in Fig. 2.5.

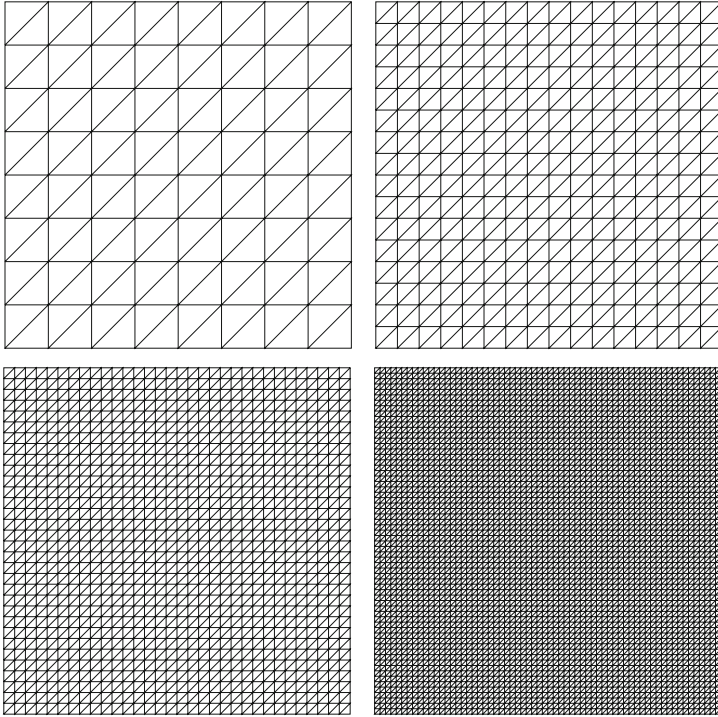


Fig. 2.4 Computational grids used for the numerical solution of problems (2.196) and (2.201)

Table 2.1 Computational errors and EOC in the $L^2(\Omega)$ -norm for the regular solution of problem (2.196)

p	$h/\sqrt{2}$	SIPG		NIPG		IIPG	
		$\ e_h\ _{L^2(\Omega)}$	EOC	$\ e_h\ _{L^2(\Omega)}$	EOC	$\ e_h\ _{L^2(\Omega)}$	EOC
1	1/8	6.7452E-02	–	2.9602E-02	–	6.3939E-02	–
1	1/16	1.8745E-02	1.85	7.6200E-03	1.96	1.7383E-02	1.88
1	1/32	4.8463E-03	1.95	1.9292E-03	1.98	4.4579E-03	1.96
1	1/64	1.2252E-03	1.98	4.8536E-04	1.99	1.1239E-03	1.99
	GEOC		1.93		1.98		1.95
2	1/8	3.9160E-03	–	1.0200E-02	–	4.7447E-03	–
2	1/16	4.9164E-04	2.99	2.5723E-03	1.99	8.4877E-04	2.48
2	1/32	6.1644E-05	3.00	6.4259E-04	2.00	1.8081E-04	2.23
2	1/64	7.7184E-06	3.00	1.6032E-04	2.00	4.2670E-05	2.08
	GEOC		3.00		2.00		2.26
3	1/8	3.1751E-04	–	5.5550E-04	–	3.2684E-04	–
3	1/16	1.9150E-05	4.05	3.4481E-05	4.01	2.0077E-05	4.02
3	1/32	1.1775E-06	4.02	2.1333E-06	4.01	1.2414E-06	4.02
3	1/64	7.3124E-08	4.01	1.3250E-07	4.01	7.7176E-08	4.01
	GEOC		4.03		4.01		4.02
4	1/8	2.3496E-05	–	3.7990E-05	–	2.7046E-05	–
4	1/16	7.5584E-07	4.96	2.4304E-06	3.97	1.2929E-06	4.39
4	1/32	2.3824E-08	4.99	1.5512E-07	3.97	7.2190E-08	4.16
4	1/64	7.4627E-10	5.00	9.7626E-09	3.99	4.3310E-09	4.06
	GEOC		4.98		3.97		4.20
5	1/8	1.4133E-06	–	2.3017E-06	–	1.6501E-06	–
5	1/16	2.2193E-08	5.99	3.6590E-08	5.98	2.6160E-08	5.98
5	1/32	3.4686E-10	6.00	5.7147E-10	6.00	4.0753E-10	6.00
5	1/64	5.4139E-12	6.00	8.8468E-12	6.01	6.3670E-12	6.00
	GEOC		6.00		6.00		6.00
6	1/8	7.3313E-08	–	1.1239E-07	–	9.5990E-08	–
6	1/16	5.8381E-10	6.97	1.5138E-09	6.21	1.1620E-09	6.37
6	1/32	4.5855E-12	6.99	2.2864E-11	6.05	1.6380E-11	6.15
6	1/64	3.8771E-14	6.89	3.5354E-13	6.02	2.4417E-13	6.07
	GEOC		6.95		6.09		6.19

We observe that EOC of the SIPG technique are in a good agreement with the theoretical ones, i.e., $O(h^{p+1})$ in the $L^2(\Omega)$ -norm (estimate (2.162)) and $O(h^p)$ in the $H^1(\Omega, \mathcal{T}_h)$ -seminorm (estimate (2.148)). On the other hand, the experimental order of convergence of the NIPG and IIPG techniques measured in the $L^2(\Omega)$ -norm is better than the theoretical estimate (2.152). We deduce that

Table 2.2 Computational errors and EOC in the $H^1(\Omega, \mathcal{T}_h)$ -seminorm for the regular solution of problem (2.196)

p	$h/\sqrt{2}$	SIPG		NIPG		IIPG	
		$ e_h _{H^1(\Omega, \mathcal{T}_h)}$	EOC	$ e_h _{H^1(\Omega, \mathcal{T}_h)}$	EOC	$ e_h _{H^1(\Omega, \mathcal{T}_h)}$	EOC
1	1/8	1.5018E+00	–	1.2423E+00	–	1.4946E+00	–
1	1/16	7.7679E–01	0.95	6.4615E–01	0.94	7.7519E–01	0.95
1	1/32	3.9214E–01	0.99	3.2741E–01	0.98	3.9181E–01	0.98
1	1/64	1.9666E–01	1.00	1.6450E–01	0.99	1.9658E–01	1.00
	GEOC		0.98		0.97		0.98
2	1/8	2.4259E–01	–	1.9985E–01	–	2.1634E–01	–
2	1/16	6.2760E–02	1.95	5.0217E–02	1.99	5.5693E–02	1.96
2	1/32	1.5849E–02	1.99	1.2536E–02	2.00	1.4053E–02	1.99
2	1/64	3.9743E–03	2.00	3.1305E–03	2.00	3.5244E–03	2.00
	GEOC		1.98		2.00		1.98
3	1/8	2.5610E–02	–	2.4029E–02	–	2.3425E–02	–
3	1/16	3.2202E–03	2.99	3.0531E–03	2.98	2.9699E–03	2.98
3	1/32	4.0238E–04	3.00	3.8298E–04	2.99	3.7253E–04	3.00
3	1/64	5.0260E–05	3.00	4.7890E–05	3.00	4.6607E–05	3.00
	GEOC		3.00		2.99		2.99
4	1/8	2.2049E–03	–	2.2096E–03	–	2.0645E–03	–
4	1/16	1.4023E–04	3.97	1.3801E–04	4.00	1.3039E–04	3.98
4	1/32	8.8035E–06	3.99	8.5962E–06	4.00	8.1650E–06	4.00
4	1/64	5.5077E–07	4.00	5.3601E–07	4.00	5.1038E–07	4.00
	GEOC		3.99		4.00		3.99
5	1/8	1.5680E–04	–	1.6457E–04	–	1.5090E–04	–
5	1/16	4.9305E–06	4.99	5.1666E–06	4.99	4.7527E–06	4.99
5	1/32	1.5413E–07	5.00	1.6126E–07	5.00	1.4865E–07	5.00
5	1/64	4.8146E–09	5.00	5.0316E–09	5.00	4.6439E–09	5.00
	GEOC		5.00		5.00		5.00
6	1/8	9.5245E–06	–	1.0198E–05	–	9.3719E–06	–
6	1/16	1.5092E–07	5.98	1.5951E–07	6.00	1.4762E–07	5.99
6	1/32	2.3666E–09	5.99	2.4862E–09	6.00	2.3083E–09	6.00
6	1/64	3.7008E–11	6.00	3.8770E–11	6.00	3.6051E–11	6.00
	GEOC		5.99		6.00		6.00

$$\|e_h\|_{L^2(\Omega)} = O(h^{\bar{p}}), \quad \bar{p} = \begin{cases} p+1 & \text{for } p \text{ odd,} \\ p & \text{for } p \text{ even.} \end{cases} \quad (2.200)$$

This interesting property of the NIPG and IIPG techniques was observed by many authors (cf. [183, 230]), but up to now a theoretical justification has been missing, see Sect. 2.9.3 for some comments. The EOC in the $H^1(\Omega, \mathcal{T}_h)$ -seminorm of NIPG and IIPG methods is in agreement with (2.148).

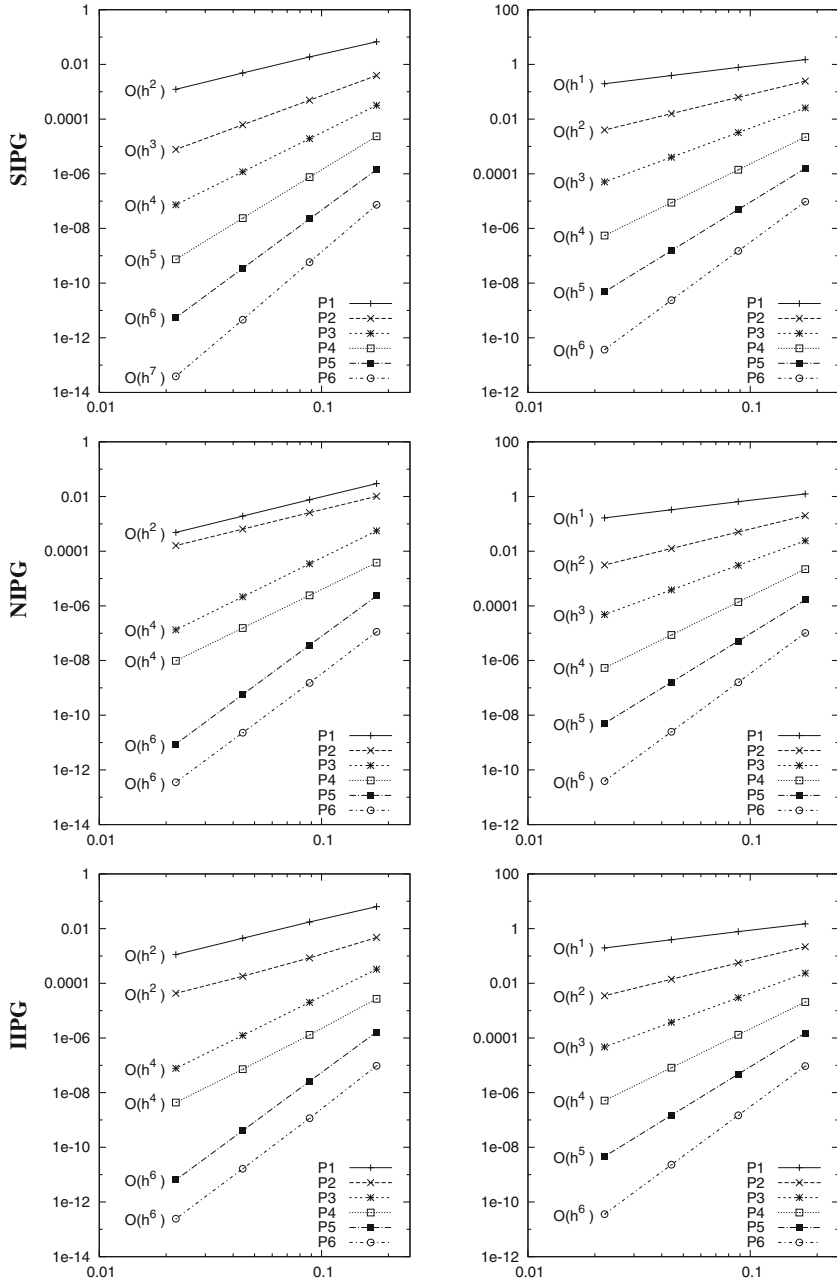


Fig. 2.5 Computational errors and EOC in the $L^2(\Omega)$ -norm (left) and in the $H^1(\Omega, \mathcal{T}_h)$ -seminorm (right) for the regular solution of problem (2.196)

2.9.2 Singular Case

In the domain $\Omega = (0, 1) \times (0, 1)$ we consider the Poisson problem

$$\begin{aligned} -\Delta u &= g \quad \text{in } \Omega, \\ u &= 0 \quad \text{on } \partial\Omega, \end{aligned} \quad (2.201)$$

with the right-hand side g chosen in such a way that the exact solution has the form

$$u(x_1, x_2) = 2r^\alpha x_1 x_2 (1 - x_1)(1 - x_2) = r^{\alpha+2} \sin(2\varphi)(1 - x_1)(1 - x_2), \quad (2.202)$$

where r, φ are the polar coordinates ($r = (x_1^2 + x_2^2)^{1/2}$) and $\alpha \in \mathbb{R}$ is a constant. The function u is equal to zero on $\partial\Omega$ and its regularity depends on the value of α . Namely, by [15],

$$u \in H^\beta(\Omega) \quad \forall \beta \in (0, \alpha + 3), \quad (2.203)$$

where $H^\beta(\Omega)$ denotes the Sobolev–Slobodetskii space of functions with *noninteger derivatives*.

We present numerical results obtained for $\alpha = -3/2$ and $\alpha = 1/2$. If $\alpha = -3/2$, then $u \in H^\beta(\Omega)$ for all $\beta \in (0, 3/2)$, whereas for the value $\alpha = 1/2$, we have $u \in H^\beta(\Omega)$ for all $\beta \in (0, 7/2)$. Figure 2.6 shows the function u for both values of α .

We carried out computations on 4 triangular grids introduced in Sect. 2.9.1 by the SIPG, NIPG and IIPG technique with the aid of P_p , $p = 1, \dots, 6$, polynomial approximations. Tables 2.3, 2.4 and Tables 2.5, 2.6 show the computational errors in the $L^2(\Omega)$ -norm as well as the $H^1(\Omega, \mathcal{T}_h)$ -seminorm, and the corresponding experimental orders of convergence for $\alpha = 1/2$ and $\alpha = -3/2$, respectively. These values are visualized in Figs. 2.7 and 2.8 in which the achieved experimental order of convergence is easy to observe.

These results lead us to the proposition that for the SIPG method the error behaves like

$$\begin{aligned} \|u - u_h\|_{L^2(\Omega)} &= O(h^\mu), & u &\in H^\beta(\Omega) \\ |u - u_h|_{H^1(\Omega)} &= O(h^{\mu-1}), & u &\in H^\beta(\Omega), \end{aligned} \quad (2.204)$$

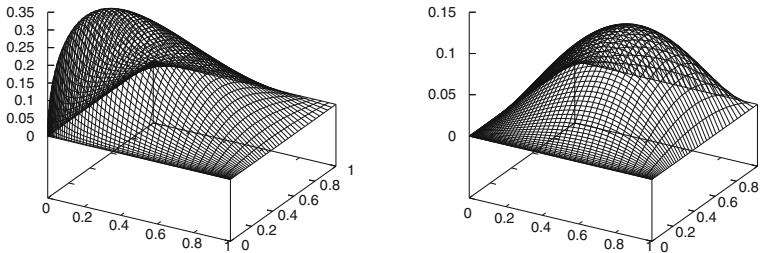


Fig. 2.6 Exact solution (2.202) for $\alpha = -3/2$ (left) and $\alpha = 1/2$ (right)

where $\mu = \min(p + 1, \beta)$, and for the IIPG and NIPG methods like

$$\begin{aligned} \|u - u_h\|_{L^2(\Omega)} &= O(h^{\bar{\mu}}), \quad u \in H^\beta(\Omega) \\ |u - u_h|_{H^1(\Omega)} &= O(h^{\mu-1}), \quad u \in H^\beta(\Omega), \end{aligned} \quad (2.205)$$

where $\mu = \min(p + 1, \beta)$, $\bar{\mu} = \min(\bar{p}, \beta)$, and \bar{p} is given by (2.200). The statements (2.204) and (2.205) are in agreement with numerical experiments (not presented here) carried out by other authors for additional values of α .

Moreover, the experimental order of convergence of the SIPG technique given by (2.204) corresponds to the result in [121], where for any $\beta \in (1, 3/2)$ we get

$$\begin{aligned} \|v - I_h v\|_{L^2(\Omega)} &\leq C(\beta) h^\mu \|v\|_{H^\beta(\Omega)}, \quad v \in H^\beta(\Omega), \\ |v - I_h v|_{H^1(\Omega)} &\leq C(\beta) h^{\mu-1} \|v\|_{H^\beta(\Omega)}, \quad v \in H^\beta(\Omega), \end{aligned} \quad (2.206)$$

where $I_h v$ is a piecewise polynomial Lagrange interpolation to v of degree $\leq p$, $\mu = \min(p + 1, \beta)$ and $C(\beta)$ is a constant independent of h and v . By [13, Sect. 3.3] and the references therein, where the interpolation in the so-called Besov spaces is used, the precise error estimate of order $O(h^{3/2})$ in the $L^2(\Omega)$ -norm and $O(h^{1/2})$ in the $H^1(\Omega, \mathcal{T}_h)$ -seminorm can be established, which corresponds to our numerical experiments.

Finally, the experimental order of convergence of the NIPG and IIPG techniques given by (2.205) corresponds to (2.206) and results (2.200).

2.9.3 A Note on the $L^2(\Omega)$ -Optimality of NIPG and IIPG

Numerical experiments from Sect. 2.9.1 lead us to the observation (2.200), which was presented, e.g., in [12, 238] and the references cited therein. The optimal order of convergence for the odd degrees of approximation was theoretically justified in [211], where NIPG and IIPG methods were analyzed for uniform partitions of the one-dimensional domain. See also [50], where similar results were obtained.

On the other hand, several examples of 1D special non-uniform (but quasi-uniform) meshes were presented in [157], where the NIPG method gives the error in the $L^2(\Omega)$ -norm of order $O(h^p)$ even for odd p . A suboptimal EOC can also be obtained for the IIPG method using these meshes, see [238, Sect. 1.5, Table 1.2].

In [101], it was shown that the use of odd degrees of polynomial approximation of IIPG method leads to the optimal order of convergence in the $L^2(\Omega)$ -norm on 1D quasi-uniform grids if and only if the penalty parameter (of order $O(h^{-1})$) is chosen in a special way. These results lead us to the hypothesis that the observation (2.200) is not valid in general.

However, extending theoretical results either to NIPG method or to higher dimensions is problematic. Some attempt was presented in [82], where the optimal order of convergence in the $L^2(\Omega)$ -norm on equilateral triangular grids was proved for the IIPG method with reduced interior and boundary penalties.

Table 2.3 Computational errors and EOC in the $L^2(\Omega)$ -norm for the solution of problem (2.201) with $\alpha = 1/2$

p	$h/\sqrt{2}$	SIPG		NIPG		IIPG	
		$\ e_h\ _{L^2(\Omega)}$	EOC	$\ e_h\ _{L^2(\Omega)}$	EOC	$\ e_h\ _{L^2(\Omega)}$	EOC
1	1/8	2.1789E-03	–	8.1338E-04	–	1.8698E-03	–
1	1/16	5.7581E-04	1.92	2.1069E-04	1.95	4.8403E-04	1.95
1	1/32	1.4740E-04	1.97	5.3806E-05	1.97	1.2267E-04	1.98
1	1/64	3.7248E-05	1.98	1.3609E-05	1.98	3.0848E-05	1.99
	GEOC		1.96		1.97		1.97
2	1/8	5.7796E-05	–	1.0098E-04	–	5.9762E-05	–
2	1/16	7.2545E-06	2.99	2.6758E-05	1.92	1.1004E-05	2.44
2	1/32	9.1150E-07	2.99	6.9525E-06	1.94	2.4341E-06	2.18
2	1/64	1.1434E-07	2.99	1.7734E-06	1.97	5.8760E-07	2.05
	GEOC		2.99		1.94		2.22
3	1/8	2.6233E-06	–	4.0597E-06	–	2.7474E-06	–
3	1/16	1.9366E-07	3.76	3.3583E-07	3.60	2.1985E-07	3.64
3	1/32	1.4898E-08	3.70	2.8012E-08	3.58	1.7889E-08	3.62
3	1/64	1.1930E-09	3.64	2.3717E-09	3.56	1.4838E-09	3.59
	GEOC		3.70		3.58		3.62
4	1/8	2.6498E-07	–	4.1937E-07	–	3.0663E-07	–
4	1/16	2.1097E-08	3.65	3.4292E-08	3.61	2.4522E-08	3.64
4	1/32	1.7819E-09	3.57	2.8705E-09	3.58	2.0460E-09	3.58
4	1/64	1.5429E-10	3.53	2.4482E-10	3.55	1.7516E-10	3.55
	GEOC		3.58		3.58		3.59
5	1/8	5.8491E-08	–	9.3494E-08	–	7.2011E-08	–
5	1/16	4.9611E-09	3.56	8.1022E-09	3.53	6.1832E-09	3.54
5	1/32	4.2999E-10	3.53	7.0989E-10	3.51	5.3944E-10	3.52
5	1/64	3.7656E-11	3.51	6.2465E-11	3.51	4.7387E-11	3.51
	GEOC		3.53		3.52		3.52
6	1/8	1.9318E-08	–	2.9767E-08	–	2.6495E-08	–
6	1/16	1.6677E-09	3.53	2.6000E-09	3.52	2.3079E-09	3.52
6	1/32	1.4570E-10	3.52	2.2856E-10	3.51	2.0259E-10	3.51
6	1/64	1.2809E-11	3.51	2.0149E-11	3.50	1.7847E-11	3.50
	GEOC		3.52		3.51		3.51

Table 2.4 Computational errors and EOC in the $H^1(\Omega, \mathcal{T}_h)$ -seminorm for the solution of problem (2.201) with $\alpha = 1/2$

p	$h/\sqrt{2}$	SIPG		NIPG		IIPG	
		$ e_h _{H^1(\Omega, \mathcal{T}_h)}$	EOC	$ e_h _{H^1(\Omega, \mathcal{T}_h)}$	EOC	$ e_h _{H^1(\Omega, \mathcal{T}_h)}$	EOC
1	1/8	5.0805E-02	–	4.2283E-02	–	5.0531E-02	–
1	1/16	2.5722E-02	0.98	2.1564E-02	0.97	2.5653E-02	0.98
1	1/32	1.2919E-02	0.99	1.0877E-02	0.99	1.2902E-02	0.99
1	1/64	6.4715E-03	1.00	5.4607E-03	0.99	6.4674E-03	1.00
	GEOC		0.99		0.98		0.99
2	1/8	4.0313E-03	–	3.2281E-03	–	3.5738E-03	–
2	1/16	1.0230E-03	1.98	8.0878E-04	2.00	9.0960E-04	1.97
2	1/32	2.5750E-04	1.99	2.0223E-04	2.00	2.2938E-04	1.99
2	1/64	6.4585E-05	2.00	5.0547E-05	2.00	5.7592E-05	1.99
	GEOC		1.99		2.00		1.99
3	1/8	2.2371E-04	–	2.2267E-04	–	2.0664E-04	–
3	1/16	3.2897E-05	2.77	3.2455E-05	2.78	3.0237E-05	2.77
3	1/32	5.0341E-06	2.71	4.9281E-06	2.72	4.5992E-06	2.72
3	1/64	8.0276E-07	2.65	7.8150E-07	2.66	7.2933E-07	2.66
	GEOC		2.71		2.72		2.72
4	1/8	2.8019E-05	–	2.6863E-05	–	2.3759E-05	–
4	1/16	4.5630E-06	2.62	4.3388E-06	2.63	3.8426E-06	2.63
4	1/32	7.7950E-07	2.55	7.3892E-07	2.55	6.5504E-07	2.55
4	1/64	1.3572E-07	2.52	1.2850E-07	2.52	1.1398E-07	2.52
	GEOC		2.56		2.57		2.57
5	1/8	8.0765E-06	–	8.3686E-06	–	7.0904E-06	–
5	1/16	1.3891E-06	2.54	1.4415E-06	2.54	1.2239E-06	2.53
5	1/32	2.4249E-07	2.52	2.5191E-07	2.52	2.1413E-07	2.51
5	1/64	4.2611E-08	2.51	4.4293E-08	2.51	3.7673E-08	2.51
	GEOC		2.52		2.52		2.52
6	1/8	3.2423E-06	–	3.4916E-06	–	2.9734E-06	–
6	1/16	5.6456E-07	2.52	6.0843E-07	2.52	5.1885E-07	2.52
6	1/32	9.9090E-08	2.51	1.0684E-07	2.51	9.1177E-08	2.51
6	1/64	1.7456E-08	2.50	1.8826E-08	2.50	1.6072E-08	2.50
	GEOC		2.51		2.51		2.51

Table 2.5 Computational errors and EOC in the $L^2(\Omega)$ -norm for the solution of problem (2.201) with $\alpha = -3/2$

p	$h/\sqrt{2}$	SIPG		NIPG		IIPG	
		$\ e_h\ _{L^2(\Omega)}$	EOC	$\ e_h\ _{L^2(\Omega)}$	EOC	$\ e_h\ _{L^2(\Omega)}$	EOC
1	1/8	9.2233E-03	–	1.4850E-02	–	7.9896E-03	–
1	1/16	3.2898E-03	1.49	5.3458E-03	1.47	2.8145E-03	1.51
1	1/32	1.1569E-03	1.51	1.8699E-03	1.52	9.8230E-04	1.52
1	1/64	4.0594E-04	1.51	6.5039E-04	1.52	3.4327E-04	1.52
	GEOC		1.50		1.51		1.51
2	1/8	2.3410E-03	–	4.6812E-03	–	1.7779E-03	–
2	1/16	8.1979E-04	1.51	1.6138E-03	1.54	6.0110E-04	1.56
2	1/32	2.8885E-04	1.50	5.6696E-04	1.51	2.0820E-04	1.53
2	1/64	1.0199E-04	1.50	2.0059E-04	1.50	7.2989E-05	1.51
	GEOC		1.51		1.51		1.53
3	1/8	9.7871E-04	–	3.1394E-03	–	1.0279E-03	–
3	1/16	3.4597E-04	1.50	1.1136E-03	1.50	3.6119E-04	1.51
3	1/32	1.2235E-04	1.50	3.9426E-04	1.50	1.2736E-04	1.50
3	1/64	4.3269E-05	1.50	1.3948E-04	1.50	4.4971E-05	1.50
	GEOC		1.50		1.50		1.50
4	1/8	6.4002E-04	–	1.6788E-03	–	7.8547E-04	–
4	1/16	2.2608E-04	1.50	5.9262E-04	1.50	2.7649E-04	1.51
4	1/32	7.9902E-05	1.50	2.0934E-04	1.50	9.7529E-05	1.50
4	1/64	2.8245E-05	1.50	7.3980E-05	1.50	3.4442E-05	1.50
	GEOC		1.50		1.50		1.50
5	1/8	3.8770E-04	–	1.1048E-03	–	6.0190E-04	–
5	1/16	1.3695E-04	1.50	3.9046E-04	1.50	2.1214E-04	1.50
5	1/32	4.8400E-05	1.50	1.3801E-04	1.50	7.4886E-05	1.50
5	1/64	1.7109E-05	1.50	4.8784E-05	1.50	2.6455E-05	1.50
	GEOC		1.50		1.50		1.50
6	1/8	2.7881E-04	–	7.5211E-04	–	5.2298E-04	–
6	1/16	9.8519E-05	1.50	2.6580E-04	1.50	1.8457E-04	1.50
6	1/32	3.4822E-05	1.50	9.3954E-05	1.50	6.5195E-05	1.50
6	1/64	1.2310E-05	1.50	3.3215E-05	1.50	2.3039E-05	1.50
	GEOC		1.50		1.50		1.50

Table 2.6 Computational errors and EOC in the $H^1(\Omega, \mathcal{T}_h)$ -seminorm for the solution of problem (2.201) with $\alpha = -3/2$

p	$h/\sqrt{2}$	SIPG		NIPG		IIPG	
		$ e_h _{H^1(\Omega, \mathcal{T}_h)}$	EOC	$ e_h _{H^1(\Omega, \mathcal{T}_h)}$	EOC	$ e_h _{H^1(\Omega, \mathcal{T}_h)}$	EOC
1	1/8	4.0604E-01	–	3.9606E-01	–	4.0035E-01	–
1	1/16	2.8999E-01	0.49	2.8508E-01	0.47	2.8631E-01	0.48
1	1/32	2.0555E-01	0.50	2.0312E-01	0.49	2.0309E-01	0.50
1	1/64	1.4539E-01	0.50	1.4413E-01	0.50	1.4370E-01	0.50
	GEOC		0.49		0.49		0.49
2	1/8	1.9294E-01	–	2.3736E-01	–	1.8460E-01	–
2	1/16	1.3627E-01	0.50	1.6750E-01	0.50	1.3052E-01	0.50
2	1/32	9.6419E-02	0.50	1.1842E-01	0.50	9.2389E-02	0.50
2	1/64	6.8224E-02	0.50	8.3741E-02	0.50	6.5385E-02	0.50
	GEOC		0.50		0.50		0.50
3	1/8	1.4304E-01	–	2.3656E-01	–	1.5217E-01	–
3	1/16	1.0145E-01	0.50	1.6731E-01	0.50	1.0794E-01	0.50
3	1/32	7.1853E-02	0.50	1.1833E-01	0.50	7.6459E-02	0.50
3	1/64	5.0852E-02	0.50	8.3679E-02	0.50	5.4113E-02	0.50
	GEOC		0.50		0.50		0.50
4	1/8	9.4937E-02	–	1.7438E-01	–	1.0791E-01	–
4	1/16	6.7297E-02	0.50	1.2334E-01	0.50	7.6474E-02	0.50
4	1/32	4.7649E-02	0.50	8.7229E-02	0.50	5.4139E-02	0.50
4	1/64	3.3715E-02	0.50	6.1686E-02	0.50	3.8306E-02	0.50
	GEOC		0.50		0.50		0.50
5	1/8	7.8490E-02	–	1.4046E-01	–	9.6583E-02	–
5	1/16	5.5605E-02	0.50	9.9348E-02	0.50	6.8396E-02	0.50
5	1/32	3.9357E-02	0.50	7.0261E-02	0.50	4.8400E-02	0.50
5	1/64	2.7843E-02	0.50	4.9686E-02	0.50	3.4238E-02	0.50
	GEOC		0.50		0.50		0.50
6	1/8	6.4288E-02	–	1.2563E-01	–	9.3368E-02	–
6	1/16	4.5518E-02	0.50	8.8855E-02	0.50	6.6077E-02	0.50
6	1/32	3.2208E-02	0.50	6.2836E-02	0.50	4.6744E-02	0.50
6	1/64	2.2782E-02	0.50	4.4434E-02	0.50	3.3060E-02	0.50
	GEOC		0.50		0.50		0.50

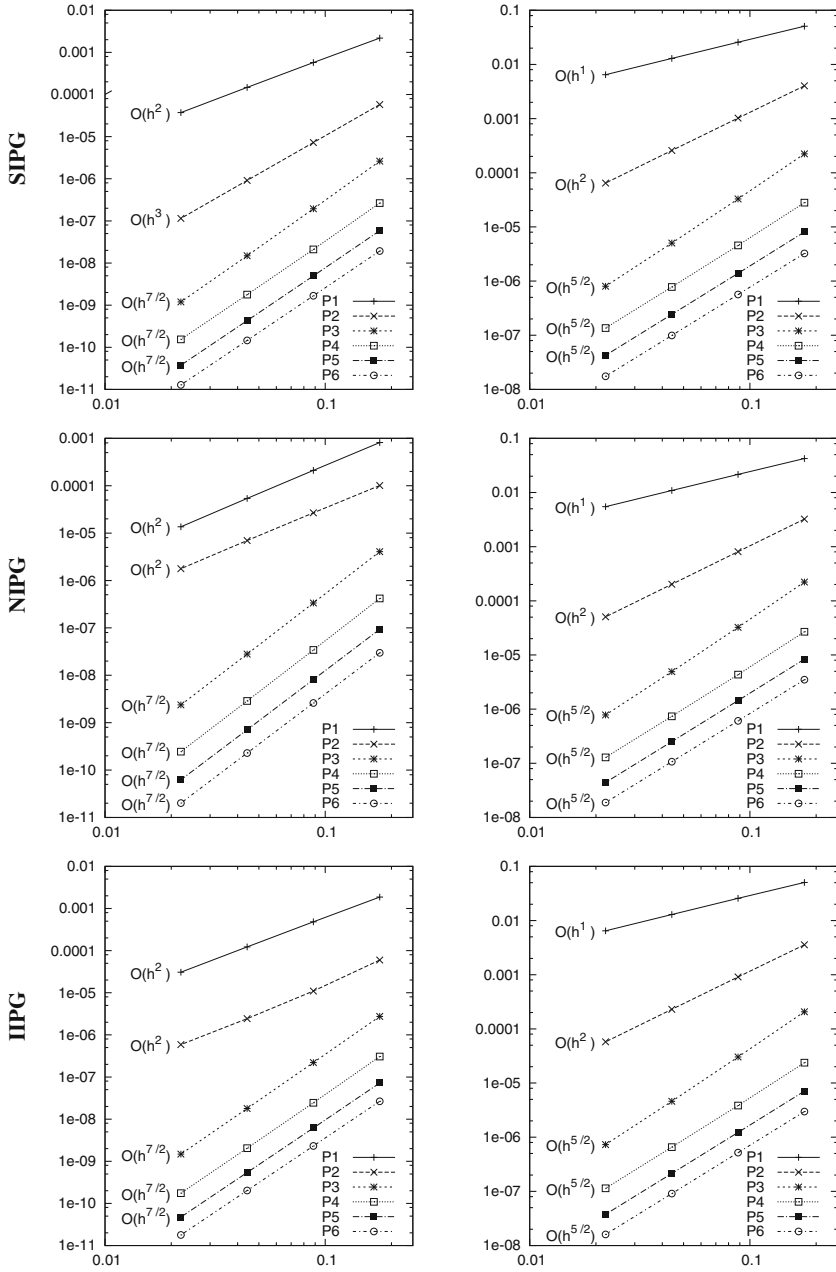


Fig. 2.7 Computational errors and EOC in the $L^2(\Omega)$ -norm (left) and the $H^1(\Omega, \mathcal{T}_h)$ -seminorm (right) for the the solution of problem (2.201) with $\alpha = 1/2$

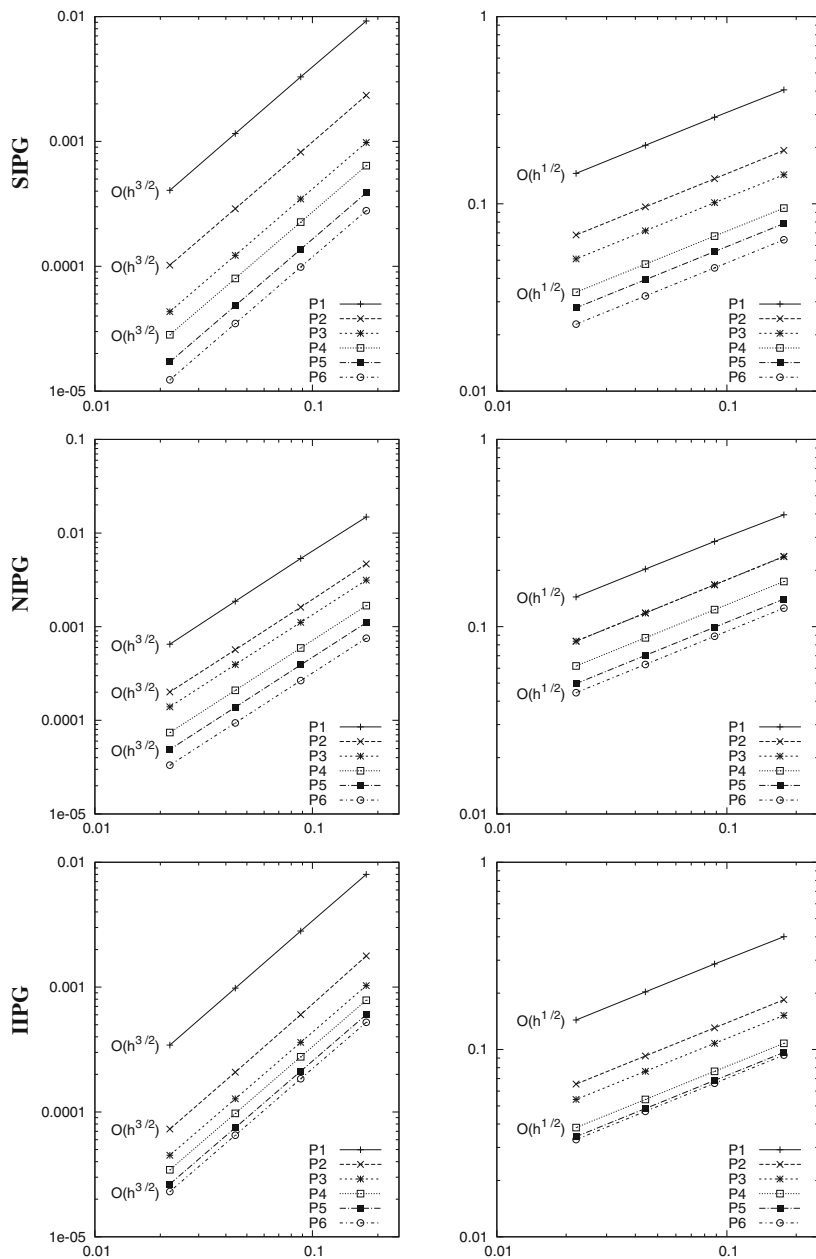


Fig. 2.8 Computational errors and EOC in the $L^2(\Omega)$ -norm (left) and the $H^1(\Omega, \mathcal{T}_h)$ -seminorm (right) for the the solution of problem (2.201) with $\alpha = -3/2$

Discontinuous Galerkin Method

Analysis and Applications to Compressible Flow

Dolejší, V.; Feistauer, M.

2015, XIV, 572 p. 87 illus., 4 illus. in color., Hardcover

ISBN: 978-3-319-19266-6