ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

# High Order Discontinuous Galerkin Method

Semester Project
of
Benjamin Stamm

Directed by:
Prof. A. Quarteroni
Dr. E. Burman

Section of Mathematics, EPFL, Lausanne

June 24, 2004

# Contents

# 1 Motivation

Standard continuous Galerkin-based finite element methods have poor stability properties when applied to transport-dominated flow problems, so excessive numerical stabilization is needed. In contrast, the Discontinuous Galerkin method is known to have good stability properties when applied to first order hyperbolic problems.

# 2 Outline

In this semester project we will consider the Discontinuous Galerkin method.

In section 3, the transport-reaction problem is presented. It mentions also the hypothesis under which uniqueness and existence of the variational formulation in $L^2$ is guaranteed. A stability result follows.

In section 4, the pure transport problem is presented as well as a stability result for this problem.

Section 5 introduces the Discontinuous Galerkin method followed by a convergence analysis. The main result of this project is the proof of the convergence theorem.

Section 6 deals with numerical results. In the framework of this semester project a Matlab code is developed for the computation of a numerical approximation. The method is applied to a simple test case with known solution. Finally, the results are analysed.

Section 7 is the conclusion of the Discontinuous Galerkin method.

In sections 8 and 9 we give a rudimentary introduction to orthogonal polynomials and numerical integration.

# 3 The Transport-Reaction Problem

In this section, the transport-reaction problem is studied with non constant coefficients. The following problem is considered:

Find $u : \Omega \to \mathbb{R}$ such that:

$$
\begin{aligned}
u_\beta + \mu u &= f &&\text{in} \quad \Omega &&(1)\\
u &= g &&\text{on} \quad \Gamma_- &&(2)
\end{aligned}
$$

where $u_\beta = \beta \cdot \nabla u$ denotes the derivative in the $\beta$-direction. $\Gamma_-$ is defined by $\Gamma_- = \{x \in \partial\Omega : n(x) \cdot \beta < 0\}$, where $n(x)$ is the outward normal unit vector at the point $x$. Analogously $\Gamma_+$ is defined by $\Gamma_+ = \{x \in \partial\Omega : n(x) \cdot \beta \geq 0\}$. $\beta$ is a vector field such that $\beta \in [W^{1,\infty}(\Omega)]^d$, and $\mu \in L^\infty(\Omega)$, $f \in H^1(\Omega)$, $g \in H^{\frac{1}{2}}(\Gamma_-)$. Let be

$$
W_\Omega = \{w \in L^2(\Omega) : \beta \cdot \nabla w \in L^2(\Omega)\} \subset L^2(\Omega)
$$

with norm

$$
\|u\|^2_{W,\Omega} = \|u\|^2_{L^2(\Omega)} + \|\beta \cdot \nabla u\|^2_{L^2\Omega}
$$

$W$ is a Hilbert space.

## 3.1   Uniqueness and existence of the variational formulation in $L^2(\Omega)$

Multiplying (1) by a smooth test function $v : \Omega \to \mathbb{R}$ and integrating on the domain $\Omega$ leads to:

Find $u \in W(K)$

$$\int_\Omega (\mu u + \beta \cdot \nabla u)v = \int_\Omega fv \qquad \forall\, v \in L^2(\Omega)$$

Let $V$ be the space

$$V = \{w \in W_\Omega : w|_{\Gamma_-} = g\} \subset W_\Omega \tag{3}$$

The bilinear form $a : W_\Omega \times L^2(\Omega) \to \mathbb{R}$ and the linear form $F : L^2(\Omega) \to \mathbb{R}$ are defined by

$$a(u,v) \;=\; \int_\Omega (\mu u + \beta \cdot \nabla u)v \qquad \forall u \in W_\Omega, \forall v \in L^2(\Omega)$$

$$F(v) \;=\; \int_\Omega vf \qquad \forall v \in L^2(\Omega)$$

Then, the variational formulation in $L^2(\Omega)$ is:

Find $u \in V$ such that:

$$a(u,v) = F(v) \qquad \forall v \in L^2(\Omega) \tag{4}$$

It can be shown that under the hypothesis that there exists a constant $\mu_0$ such that

$$\mu(x) - \frac{1}{2}\nabla \cdot \beta(x) \geq \mu_0 > 0 \quad \text{a.e. in} \quad \Omega \tag{5}$$

the conditions of the Nečas Theorem are satisfied. That implies that there exists a unique solution of (4). In addition, condition (5) is also necessary for uniqueness and existence. For more details see [1].

## 3.2   Stability

In this section, a stability result for the Transport-Reaction Problem on the whole domain is developed.

**Lemma 3.1 (stability for the transport-reaction problem)** *If it exists a constant $\mu_1$ such that $\mu(x) \geq \mu_1 > 0 \;\forall\, x \in \Omega$, then the following stability result is given:*

$$\mu_1 \|u\|^2_{L^2(\Omega)} + \int_{\Gamma_+} |\beta \cdot n| u^2 \leq \frac{1}{\mu_1}\|f\|^2_{L^2(\Omega)} + \int_{\Gamma_-} |\beta \cdot n| g^2$$

PROOF. Let us take equation (1), multiply it by $u$ and integrate over $\Omega$. We get

$$(u_\beta, u)_\Omega + (\mu u, u)_\Omega = (f, u)_\Omega$$

Then, the following integration by parts is used

$$\begin{aligned} (u_\beta, u)_\Omega \;&=\; -(u, u_\beta)_\Omega + (\beta \cdot nu, u)_{\partial\Omega} \\ &=\; -(u_\beta, u)_\Omega + (|\beta \cdot n|u, u)_{\Gamma_+} - (|\beta \cdot n|g, g)_{\Gamma_-} \end{aligned}$$

so that

$$(\mu u, u)_\Omega + \frac{1}{2}(|\beta \cdot n|u, u)_{\Gamma_+} = (f, u)_\Omega + \frac{1}{2}(|\beta \cdot n|g, g)_{\Gamma_-}$$

Now, the Cauchy-Schwarz inequality is applied to $(f, u)_\Omega$. This leads to

$$2(\mu u, u)_\Omega + (|\beta \cdot n|u, u)_{\Gamma_+} = 2\|f\|_\Omega \|f\|_\Omega + (|\beta \cdot n|g, g)_{\Gamma_-}$$

Finally the Young inequality with $\varepsilon = \frac{1}{2\mu_1}$ is applied. Then

$$\mu_1 \|u\|_{L^2(\Omega)}^2 + \int_{\Gamma_+} |\beta \cdot n||u|^2 \leq \frac{1}{\mu_1}\|f\|_{L^2(\Omega)}^2 + \int_{\Gamma_-} |\beta \cdot n||g|^2$$

$\square$

# 4 The pure Transport Problem

In this section, the pure Transport Problem is considered:
Find $u : \Omega \to \mathbb{R}$ such that

$$
\begin{align}
u_\beta &= f \quad \text{in} \quad \Omega \subset \mathbb{R}^d \tag{6} \\
u &= g \quad \text{on} \quad \Gamma_- \tag{7}
\end{align}
$$

Then the variational formulation of this problem is:
Find $u \in W_\Omega$ such that:

$$
\begin{align}
(u_\beta, v)_\Omega &= (f, v)_\Omega \quad \forall v \in L^2(\Omega) \tag{8} \\
u &= g \quad \text{on} \quad \Gamma_-
\end{align}
$$

## 4.1 Stability

**Lemma 4.1 (stability for the pure Transport Problem)** *In the case of the pure transport problem, i.e. if $\mu \equiv 0$ and under the hypothesis that there exists a vector function $\gamma \in \left(L^\infty(K)\right)^d$ such that it exists a constant $\gamma_1$ which satisfies*

$$\vec{\gamma} \cdot \vec{\beta} \geq \gamma_1 > 0$$

*then, the following stability result is given*

$$\gamma_1 \|e^{-\vec{\gamma} \cdot \vec{x}} u\|_{L^2(\Omega)}^2 + (|\beta \cdot n|e^{-\vec{\gamma} \cdot \vec{x}} u, e^{-\vec{\gamma} \cdot \vec{x}} u)_{\Gamma_+}$$
$$\leq \frac{1}{\gamma_1}\|e^{-\vec{\gamma} \cdot \vec{x}} f\|_{L^2(\Omega)}^2 + (|\beta \cdot n|e^{-\vec{\gamma} \cdot \vec{x}} g, e^{-\vec{\gamma} \cdot \vec{x}} g)_{\Gamma_-}$$

PROOF. Let be $\tilde{u}(\vec{x}) = e^{-\vec{\gamma} \cdot \vec{x}} u(\vec{x})$ and note that problem (6) is equivalent to the following problem:
Find $\tilde{u} : \Omega \to \mathbb{R}$ such that

$$
\begin{align}
\vec{\beta} \cdot \nabla \tilde{u} + \nu \tilde{u} &= e^{-\vec{\gamma} \cdot \vec{x}} f \quad \text{in} \quad \Omega \\
\tilde{u} &= e^{-\vec{\gamma} \cdot \vec{x}} g \quad \text{on} \quad \Gamma_-
\end{align}
$$

where $\nu = \vec{\gamma} \cdot \vec{\beta}$. This is a transport-reaction problem and thanks to the hypothesis, Lemma (3.1) can be applied, so that we get the result.

$\square$

# 5 Discontinuous Galerkin Method

## 5.1 Notations

Let us first introduce some notations. Consider an element $K$. $K$ can be arbitrary a simplex or a parallelepiped , then $\partial K$ is split into

$$
\begin{aligned}
\partial K_- &= \{x \in \partial K : n(x) \cdot \beta < 0\} \\
\partial K_+ &= \{x \in \partial K : n(x) \cdot \beta \geq 0\}
\end{aligned}
$$

and we have that $\partial K = \partial K_- \cup \partial K_+$.

$$
\begin{aligned}
v^- &= \lim_{s \to 0^-} v(x + s\beta) \\
v^+ &= \lim_{s \to 0^+} v(x + s\beta)
\end{aligned}
$$

$$
[v] = v^+ - v^-
$$

$$
\hat{v}(x) = \begin{cases} v^+ & x \in \partial K_- \\ v^- & x \in \partial K_+ \end{cases}
$$

Let be $h_K$ the diameter of the element $K$ and $h = \max_K h_K$.
In addition, $(u, v)_\Lambda$ denotes the usual $L^2$-scalar product on $\Lambda$.

## 5.2 The method

In this section, the principle of the Discontinuous Galerkin method is presented. Let us consider the transport-reaction equations. We will begin by studying the problem restricted to one element, meaning a simplex or parallelepiped. On the element $K$, the following problem is considered:

Find $u : \Omega \to \mathbb{R}$ such that

$$
\begin{aligned}
u_\beta + \mu u &= f && \text{in} \quad K & (9) \\
u &= u^- && \text{on} \quad \partial K_- & (10)
\end{aligned}
$$

where $u^-$ is a given function on $\partial K_-$. The associated variational formulation is:

Find $u \in W_K$

$$
\begin{aligned}
(u_\beta + \mu u, v)_K &= (f, v)_K && \forall v \in L^2(K) \\
u &= u^- && \text{on} \quad \partial K_-
\end{aligned}
$$

Using integration by parts leads to:

$$
(u_\beta, v)_K + (\mu u, v)_K = -(u, v_\beta)_K + (\beta \cdot n \hat{u}, \hat{v})_{\partial K} + (\mu u, v)_K = (f, v)_K
$$

where $n$ is the outward unit vector of $K$. Boundary conditions are imposed weakly, so that the problem is:

Find $u \in W_K$ such that:

$$
(\mu u, v)_K - (u, v_\beta)_K + (\beta \cdot n u^-, v^-)_{\partial K_+} = (f, v)_K + (|\beta \cdot n| u^-, v^+)_{\partial K_-} \qquad \forall v \in W_K
$$

Note that $u^-$ on $\partial K_-$ is the given function of the boundary condition, but that $u^-$ on $\partial K_+$ is the solution on $\partial K_+$.

Now, a Galerkin approximation is used. That means that the space $W_K$ is replaced by the finite dimensional space $V_h$ which satisfies $V_h \subset W_K$. The formulation is the following:

Find $u_h \in V_h$ such that:

$$(\mu u_h, v_h)_K - (u_h, v_{h\beta})_K + (\beta \cdot n u_h^-, v_h^-)_{\partial K_+}$$
$$= (f, v_h)_K + (|\beta \cdot n| u_h^-, v_h^+)_{\partial K_-} \qquad \forall v \in V_h$$

where $u_h^-$ is an approximation of $u^-$ in the space $V_h$. Then using integration by parts for $(u_h, v_{h\beta})_K$ a second time leads to:

Find $u_h \in V_h$) such that:

$$a_K(u_h, v_h) = (f, v_h)_K \qquad \forall v \in V_h \tag{11}$$

where $u_h^-$ is given on the inflow boundary $\partial K_-$ and

$$a_K(u_h, v_h) = (u_{h\beta} + \mu u_h, v_h)_K + (|\beta \cdot n|[u_h], v_h^+)_{\partial K_-}$$

Choosing a basis $\{\varphi_i\}$ for the space $V_h$ and writing $u_h = \sum_j u_j \varphi_j(x)$ leads to a algebraic system:

$$A\vec{u} = \vec{f} \tag{12}$$

where

$$A_{i,j} = (\varphi_{j,\beta} + \mu\varphi_j, \varphi_i)_K + (|\beta \cdot n|\varphi_j^+, \varphi_i^+)_{\partial K_-}$$
$$\vec{f_i} = (f, \varphi_i)_K + (|\beta \cdot n|u^-, \varphi_i^+)_{\partial K_-}$$

To solve the problem on one element, one needs only to know the the inflow boundary data and $f$. Because the data is given on the inflow boundary of the whole domain, one can find an order of elements such that the problem can be solved element by element. For an example see Fig. 1. Instead of solving a big system of linear equations, many small systems of linear equations has to be solved when using Discontinuous Galerkin.

Now, we would like to show existence and uniqueness of the approximation $u_h$. Since the problem can be uncoupled, it is sufficient to prove well posedness for one element $K$. The following lemma shows the result in the case of the Transport-Reaction problem.

**Lemma 5.1** *If $\mu \in L^\infty(\Omega)$ is such that*

$$\mu(x) \geq \mu_1 > 0 \qquad \forall\, x \in \Omega$$

*then for each element $K \in \tau_h$, it exists a unique solution of the linear system (12)*

$$A\vec{u} = \vec{f}$$

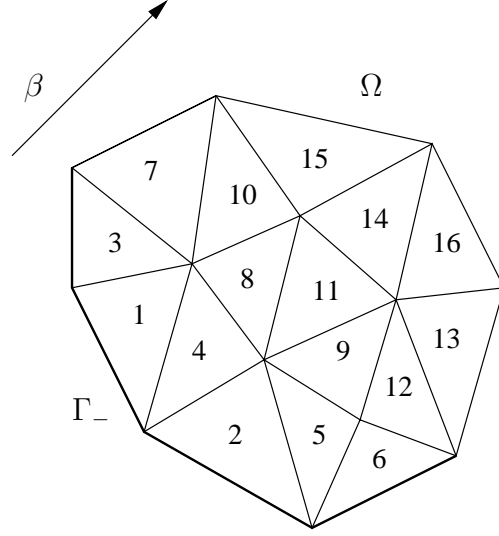*for every constant vector $\vec{f}$.*

Figure 1: Example of an order of triangles for the Discontinuous Galerkin Method

PROOF. First, the test function $v_h \in V_h$ of the variational formulation is chosen equal to $u_h$. Using partial integration, we obtain

$$(u_{h\beta}, u_h)_K = -(u_h, u_{h\beta})_K + (\beta \cdot n u_h^+, u_h^+)_{\partial K_-} + (\beta \cdot n u_h^-, u_h^-)_{\partial K_+} \tag{13}$$

so that

$$2\,(u_{h\beta}, u_h)_K = (\beta \cdot n u_h^+, u_h^+)_{\partial K^-} + (\beta \cdot n u_h^-, u_h^-)_{\partial K^+} \tag{14}$$

Secondly, the left hand side of the variational formulation is considered.

$$\begin{aligned}
2\,(u_{h\beta}, u_h)_K &+ 2\,(|\beta \cdot n| u_h^+, u_h^+)_{\partial K_-} \tag{15}\\
&= -(|\beta \cdot n| u_h^+, u_h^+)_{\partial K_-} + (|\beta \cdot n| u_h^-, u_h^-)_{\partial K_+}\\
&\quad + 2\,(|\beta \cdot n| u_h^+, u_h^+)_{\partial K_-}\\
&= (|\beta \cdot n| u_h^-, u_h^-)_{\partial K_+} + (|\beta \cdot n| u_h^+, u_h^+)_{\partial K_-}\\
&= \int_{\partial K} |\beta \cdot n| \hat{u_h}^2 d\sigma
\end{aligned}$$

so that

$$\begin{aligned}
2\,(u_{h\beta} + \mu u_h, u_h)_K &+ 2\,(|\beta \cdot n| u_h^+, u_h^+)_{\partial K_-}\\
&= \int_{\partial K} |\beta \cdot n| \hat{u_h}^2 d\sigma + 2 \int_K \mu u_h^2 dx
\end{aligned}$$

Using the same basis of the space $V_h$ as above, this is equivalent to the following equation:

$$\vec{u} A \vec{u} = \frac{1}{2} \int_{\partial K} |\beta \cdot n| u_h^2 d\sigma + \int_K \mu \hat{u_h}^2 dx \tag{16}$$

If $A\vec{u} = 0$ then $\vec{u} A \vec{u} = 0$ and as consequence

$$0 \geq \frac{1}{2} \int_{\partial K} |\beta \cdot n| \hat{u_h}^2 d\sigma + \mu_1 \int_K u_h^2 dx \geq \|u_h\|_{L^2(K)}^2 \geq 0 \tag{17}$$

That implies that $u_h \equiv 0$ and $\vec{u} = \vec{0}$. This proves the regularity of the matrix $A$ as well as the uniqueness and existence of the approximation $u_h$.

$\square$

In the case of the pure Transport Problem, we have the following result.

**Lemma 5.2** *We assume that one of the following conditions are satisfied.*

- $V_h = P^1(K)$

- $V_h = P^2(K)$ *and each side $\Gamma$ of the the triangle K satisfies*

$$|\beta \cdot n(x)| > C \qquad \forall x \in \Gamma$$

*then, there exists a unique solution of the linear system (12), derived from the pure transport problem,*
$$A\vec{u} = \vec{f}$$
*for every constant vector $\vec{f}$.*

PROOF. Taking equation (16) with $\mu = 0$ leads to:

$$\vec{u}A\vec{u} = \frac{1}{2} \int_{\partial K} |\beta \cdot n| \hat{u_h}^2 d\sigma \tag{18}$$

The uniqueness and existence of $u_h$ follows of the regularity of the matrix $A$. Let be $\vec{u}$ so that $A\vec{u} = \vec{0}$, then

$$\int_{\partial K} |\beta \cdot n| \hat{u_h}^2 d\sigma = 0 \tag{19}$$

If it can be shown that $\vec{u} = \vec{0}$, resp. $u_h = 0$, A is a regular matrix and $u_h$ exists and is unique for each $\vec{f}$. But concluding that $u_h = 0$ depends now from the space $V_h$.

- If $V_h = P^1(K)$, then $u_h = 0$, also if $\beta$ is parallel to one side of the element $K$.

- If $V_h = P^2(K)$ it is necessary that each side $\Gamma$ of the triangle K satisfies

$$|\beta \cdot n(x)| > C \qquad \forall x \in \Gamma \tag{20}$$

where $C > 0$. This is a supplementary condition of the triangulation.

$\square$

**Remark**: Figure 2 shows the interpolation points for the different spaces.
**Remark**: If $V_h = P^N(K)$ with $N \geq 3$, there is no control over the basis functions with support in the interior of the triangle. For example suppose that $N = 3$, then one interpolation point $p_{int}$ is in the interior of the triangle. Let $\psi$ be the function such that

$$\psi(p_{int}) = 1 \qquad \text{and} \qquad \psi(x) = 0 \quad \forall x \in \partial K$$

Then

$$\int_{\partial K} |\beta \cdot n| \psi^2 = 0 \quad \text{but} \quad \psi(x) \neq 0 \, \forall x \in int(K)$$

Figure 2: The interpolation points for $P^1(K)$, $P^2(K)$ and $P^3(K)$

## 5.3 Convergence Analysis

In this section, a global convergence result for the Discontinuous Galerkin method is developed under the hypothesis that $\mu = 1$. Note that in this section $C$ and $C_K$ denotes generic constants taking different values. We will not pay attention to the explicit form ot these constants but note that they are independent of $h$ and $N$, unless otherwise noted. We start from the local problem defined by (11). Then, the sum over all elements $K \in \tau_h$ is taken and. This describes the global problem:

Find $u_h \in W_h^N$ such that

$$\sum_{K \in \tau_h} a_K(u_h, v_h) = (f, v_h)_\Omega + (|\beta \cdot n|g, v_h^+)_{\Gamma_-} \qquad \forall v_h \in W_h^N$$

where

$$W_h^N = \{w \in L^2(\Omega)| \; w_{|K} \in V_h \; , \; \forall K \in \tau_h\}$$

In fact, this is equivalent to

Find $u_h \in W_h^N$ such that

$$a(u_h, v_h) = (f, v_h)_\Omega + (|\beta \cdot n|g, v_h^+)_{\Gamma_-} \qquad \forall v_h \in W_h^N \qquad (21)$$

where

$$a(u_h, v_h) = (\mu u_h + u_{h,\beta}, v_h)_\Omega + \sum_K (|\beta \cdot n|[u_h], v_h^+)_{\partial K_- \setminus \Gamma} + (|\beta \cdot n|u_h^+, v_h^+)_{\Gamma_-}$$

Note that the result is developed for $\mu = 1$. Taking $v_h = u_h$ and integrating by parts, we obtain

$$
\begin{aligned}
a(u_h, u_h) &= \|u_h\|_{L^2(\Omega)}^2 + \sum_K (u_{h,\beta}, u_h)_K + \sum_K (|\beta \cdot n|[u_h], v_h^+)_{\partial K_- \setminus \Gamma} \\
&\quad + (|\beta \cdot n|u_h^+, v_h^+)_{\Gamma_-} \\
&= \|u_h\|_{L^2(\Omega)}^2 + \sum_K (|\beta \cdot n|[u_h], u_h^+)_{\partial K_- \setminus \Gamma} + (|\beta \cdot n|u_h^+, u_h^+)_{\Gamma_-} \\
&\quad + \sum_K \frac{1}{2}\Big( (|\beta \cdot n|u_h^-, u_h^-)_{\partial K_+} - (|\beta \cdot n|u_h^+, u_h^+)_{\partial K_-} \Big)
\end{aligned}
$$

Now, the following equality is used

$$(a - b)a = \frac{1}{2}(a^2 + (a - b)^2 - b^2)$$

so that

$$(|\beta \cdot n|[u_h], u_h^+)_{\partial K_- \backslash \Gamma} = \frac{1}{2} \int_{\partial K_- \backslash \Gamma} |\beta \cdot n|((u_h^+)^2 + [u_h]^2 - (u^-)^2)$$

then

$$
\begin{aligned}
a(u_h, u_h) &= \|u_h\|_{L^2(\Omega)}^2 + \sum_K \Big[ \frac{1}{2} \int_{\partial K_- \backslash \Gamma} |\beta \cdot n| \big( (u_h^+)^2 + [u_h]^2 - (u_h^-)^2 \big) \\
&\quad + \frac{1}{2} \int_{\partial K_- \backslash \Gamma} |\beta \cdot n|(u_h^-)^2 - \frac{1}{2} \int_{\partial K_- \backslash \Gamma} |\beta \cdot n|(u_h^+)^2 \Big] \\
&\quad + \frac{1}{2} \int_\Gamma |\beta \cdot n| \hat{u_h}^2 \\
&= \|u_h\|_{L^2(\Omega)}^2 + \sum_K \Big( \frac{1}{2} \int_{\partial K_- \backslash \Gamma} |\beta \cdot n|[u_h]^2 \Big) + \frac{1}{2} \int_\Gamma |\beta \cdot n| \hat{u_h}^2
\end{aligned}
$$

This motivates us to define the following norm

$$\||u_h\||^2 = \|u_h\|_{L^2(\Omega)}^2 + \sum_K \Big( \frac{1}{2} \int_{\partial K_- \backslash \Gamma} |\beta \cdot n|[u_h]^2 \Big) + \frac{1}{2} \int_\Gamma |\beta \cdot n| \hat{u_h}^2$$

and as consequence

$$\||u_h\||^2 = a(u_h, u_h).$$

We immediately conclude that

$$\||u_h\||^2 = (f, u_h)_\Omega + (|\beta \cdot n|g, u_h^+)_{\Gamma_-}.$$

Existence and uniqueness of the discrete solution follows. For $u \in H^1(\Omega)$ the numerical schema is consistent. This is expressed in the form of the Galerkin orthogonality in the following lemme.

**Lemma 5.3 (Galerkin orthogonality)**

$$a(u_h - u, v_h) = 0 \qquad \forall\, v_h \in W_h^N$$

*where $u$ is the exact solution and $u_h$ solution of (21).*

PROOF. First, since $u_h$ is solution of (21), it satisfies

$$a(u_h, v_h) = (f, v_h)_\Omega + (|\beta \cdot n|g, v_h^+)_{\Gamma_-} \qquad \forall v_h \in W_h^N \tag{22}$$

Secondly, with $u$ the exact solution, we have

$$
\begin{aligned}
(\mu u + u_\beta, v_h)_\Omega &= (f, v_h)_\Omega \qquad \forall v_h \in W_h^N \\
(u - g, v_h)_{\Gamma_-} &= 0 \qquad \forall v_h \in W_h^N
\end{aligned}
$$

In addition, $u \in H^1(\Omega)$ and as consequence $u \in H^1(K) \ \forall K \in \tau_h$. So that the trace is well defined and hence

$$\int_{\partial K} |\beta \cdot n|[u]v \leq C \int_{\partial K} [u]v = 0.$$

So that

$$a(u, v_h) = (f, v_h)_\Omega + (|\beta \cdot n|g, v_h^+)_{\Gamma_-} \qquad \forall v_h \in W_h^N \tag{23}$$

Taking the difference between (22) and (23) leads to

$$a(u_h - u, v_h) = 0 \qquad \forall\, v_h \in W_h^N.$$

$\square$

Now, the local projection $P_N : L^2(K) \to W_h^N$ is introduced. For all $u \in L^2(K)$, there exists $P_N u \in W_h^N \cap C^0(K)$ such that

$$\int_K (u - P_N u)v_h = 0 \quad \forall v_h \in W_h^N.$$

Let be $h_K$ the diameter of the element $K$ defined by $h_K = \max_{x,y \in K} |x - y|$. Then $h$ is defined by $h = \max_K h_K$ In addition, we refer to [2] for the following Lemma's.

**Lemma 5.4 (Local Projection Error on the boundary $\partial K$)** *Let $K \in \tau_h$ and suppose that $u \in H^k(K)$ for some integer $k \geq 1$. Then, for any integer $s, 1 \leq s \leq min(N + 1, k)$ and $N \geq 1$, we have that*

$$\|u - P_N u\|_{L^2(\partial K)} \leq C_K(d, s) \frac{h_K^{s - \frac{1}{2}}}{(N + 1)^{s - \frac{1}{2}}} |u|_{H^s(K)}$$

**Lemma 5.5 (Local Projection Error)** *Let $K \in \tau_h$ and suppose that $u \in H^k(K)$ for some integer $k \geq 1$. Then, for any integer $s, 1 \leq s \leq min(N + 1, k)$ and $N \geq 1$, we have that*

$$\|u - P_N u\|_{L^2(K)} \leq C_K \frac{h_K^s}{N^s} |u|_{H^s(K)}$$

We first state and proof two approximation lemmas.

**Lemma 5.6** *If there exists a constant $C(\beta)$ such that $|\beta \cdot n| \leq C(\beta)$, then*

$$\||\eta\|| \leq \sum_K C_K(d, s, \beta) \frac{h_K^{s - \frac{1}{2}}}{N^{s - \frac{1}{2}}} |u|_{H^s(K)}$$

PROOF Lemma(5.6)

$$\||\eta\||^2 = \|\eta\|_{L^2(\Omega)}^2 + \sum_K \frac{1}{2} \int_{\partial K_- \backslash \Gamma} |\beta \cdot n| [\eta]^2 + \sum_K \frac{1}{2} \int_{\partial K \cap \Gamma} |\beta \cdot n| \hat{\eta}^2$$

Considering each term of this sum:

- $\|\eta\|_{L^2(\Omega)}^2 \leq \sum_K \left( C \frac{h_K^s}{N^s} |u|_{H^s(K)} \right)^2$ by Lemma(5.5)

- 
$$
\begin{aligned}
\sum_K \frac{1}{2} \int_{\partial K \cap \Gamma} |\beta \cdot n| \hat{\eta}^2 &\leq \sum_K \frac{1}{2} \int_{\partial K} |\beta \cdot n| \hat{\eta}^2 \\
&\leq C(\beta) \sum_K \|\eta\|_{L^2(\partial K)}^2 \leq \sum_K \left( C_K(d, s, \beta) \frac{h_k^{s - \frac{1}{2}}}{(N + 1)^{s - \frac{1}{2}}} |u|_{H^s(K)} \right)^2
\end{aligned}
$$

- 
$$
\begin{aligned}
\sum_K \frac{1}{2} \int_{\partial K_- \backslash \Gamma} |\beta \cdot n| [\eta]^2 &\leq \sum_K \frac{1}{2} \int_{\partial K_- \backslash \Gamma} |\beta \cdot n| \left( \eta^{+2} + \eta^{-2} \right) \\
&\leq 2 \sum_K \int_{\partial K} |\beta \cdot n| \hat{\eta}^2 \leq C(\beta) \sum_K \|\eta\|_{L^2(\partial K)}^2 \\
&\leq \sum_K \left( C_K(d, s, \beta) \frac{h_k^{s - \frac{1}{2}}}{(N + 1)^{s - \frac{1}{2}}} |u|_{H^s(K)} \right)^2
\end{aligned}
$$

So that

$$
\begin{aligned}
\|\|\eta\|\| &\leq \Big( \sum_K \big( C_K(d,s,\beta) \frac{h_k^{s-\frac{1}{2}}}{(N+1)^{s-\frac{1}{2}}} |u|_{H^s(K)} \big)^2 \\
&\quad + \big( C_K(d,s,\beta) \frac{h_k^{s-\frac{1}{2}}}{N^{s-\frac{1}{2}}} |u|_{H^s(K)} \big)^2 \Big)^{\frac{1}{2}} \\
&\leq \Big( \sum_K \big( C_K(d,s,\beta) \frac{h_k^{s-\frac{1}{2}}}{N^{s-\frac{1}{2}}} |u|_{H^s(K)} \big)^2 \Big)^{\frac{1}{2}} \leq \sum_K C_K(d,s,\beta) \frac{h_k^{s-\frac{1}{2}}}{N^{s-\frac{1}{2}}} |u|_{H^s(K)}
\end{aligned}
$$

$\square$ Lemma(5.6)

**Lemma 5.7** *If there exists a constant $C(\beta)$ such that $|\beta \cdot n| \leq C(\beta)$, then*

$$
\Big( \sum_K \int_{\partial K} |\beta \cdot n| (\eta^-)^2 \Big)^{\frac{1}{2}} \leq \sum_K C_K(d,s,\beta) \frac{h_K^{s-\frac{1}{2}}}{(N+1)^{s-\frac{1}{2}}} |u|_{H^s(K)}
$$

PROOF Lemma(5.7)

$$
\begin{aligned}
\Big( \sum_K \int_{\partial K} |\beta \cdot n| (\eta^-)^2 \Big)^{\frac{1}{2}} &\leq \Big( \sum_K C \|\hat{\eta}\|^2_{L^2(\partial K)} \Big)^{\frac{1}{2}} \\
&\leq \Big( \sum_K ( C_K(d,s,\beta) \frac{h_K^{s-\frac{1}{2}}}{(N+1)^{s-\frac{1}{2}}} |u|_{H^s(K)} )^2 \Big)^{\frac{1}{2}} \\
&\leq \sum_K C_K(d,s,\beta) \frac{h_K^{s-\frac{1}{2}}}{(N+1)^{s-\frac{1}{2}}} |u|_{H^s(K)}
\end{aligned}
$$

$\square$ Lemma(5.7)

**Theorem 5.8 (Global Convergence)** *Suppose that $u \in H^k(K)$ for some integer $k \geq 1$ and there exists a constant $C(\beta)$ such that $|\beta \cdot n| \leq C(\beta)$. Then, for any integer $s, 1 \leq s \leq min(N+1, k)$ and $N \geq 1$, we have that*

$$
\|\|u - u_h\|\| \leq C(d,s,\beta) \frac{h^{s-\frac{1}{2}}}{N^{s-\frac{1}{2}}} |u|_{H^s(\Omega)}
$$

PROOF. Let be $\eta = P_N u - u$ and $\xi = P_N u - u_h$. Then

$$
\|\|u - u_h\|\| \leq \|\|u - P_N u\|\| + \|\|P_N u - u_h\|\| = \|\|\eta\|\| + \|\|\xi\|\|
$$

Let us show the following inequality

$$
\|\|\xi\|\| \leq \sum_K C_K(d,s,\beta) \frac{h_K^{s-\frac{1}{2}}}{(N+1)^{s-\frac{1}{2}}} |u|_{H^s(K)}
$$

The Galerkin Orthogonality is used for developing

$$
\begin{aligned}
\|\|\xi\|\|^2 &= a(P_N u - u_h, P_N u - u_h) + a(u_h - u, P_N u - u_h) = a(\eta, \xi) \\
&= \sum_K (\eta, \xi)_K + \sum_K (\eta_\beta, \xi)_K + \sum_K (|\beta \cdot n|[\eta], \xi^+)_{\partial K_- \backslash \Gamma} \\
&\quad + (|\beta \cdot n|\eta^+, \xi^+)_{\Gamma_-}
\end{aligned}
$$

Observe that $\xi = P_N u - u_h \in W_h^N$ and as consequence

$$(\eta, \xi)_K = 0$$

Using integration by parts of $(\eta_\beta, \xi)_K$:

$$(\eta_\beta, \xi)_K = (\eta, \xi_\beta)_K + \sum_K (|\beta \cdot n| \hat{\eta}, \hat{\xi})_{\partial K}$$

and the fact that $\xi_\beta \in W_h^N$ and in consequence that $(\eta, \xi_\beta)_K = 0$ leads to

$$
\begin{aligned}
\|\|\xi\|\|^2 &= \sum_K (|\beta \cdot n|[\eta], \xi^+)_{\partial K_- \backslash \Gamma} + (|\beta \cdot n|\eta^+, \xi^+)_{\Gamma_-} \\
&\quad + \sum_K (|\beta \cdot n|\eta^-, \xi^-)_{\partial K_+} - \sum_K (|\beta \cdot n|\eta^+, \xi^+)_{\partial K_-} \\
&= \sum_K (|\beta \cdot n|\eta^-, \xi^-)_{\partial K_+} - \sum_K (|\beta \cdot n|\eta^-, \xi^+)_{\partial K_- \backslash \Gamma} \\
&= \sum_K (|\beta \cdot n|\eta^-, \xi^- - \xi^+)_{\partial K_- \backslash \Gamma} + (|\beta \cdot n|\eta^-, \xi^-)_{\Gamma_+}
\end{aligned}
$$

First, the Cauchy-Schwarz inequality is applied to

$$(|\beta \cdot n|\eta^-, \xi^-)_{\Gamma_+} \leq \left( \int_{\Gamma_+} |\beta \cdot n|(\eta^-)^2 \right)^{\frac{1}{2}} \left( \int_{\Gamma_+} |\beta \cdot n|(\xi^-)^2 \right)^{\frac{1}{2}}$$

Observing that

$$\left( \int_{\Gamma_+} |\beta \cdot n|(\xi^-)^2 \right)^{\frac{1}{2}} \leq \|\|\xi\|\|$$

and

$$\left( \int_{\Gamma_+} |\beta \cdot n|(\eta^-)^2 \right)^{\frac{1}{2}} \leq \left( \sum_K \int_{\partial K} |\beta \cdot n|(\eta^-)^2 \right)^{\frac{1}{2}}$$

Lemma(5.7) can be applied, so that

$$
\begin{aligned}
(|\beta \cdot n|\eta^-, \xi^-)_{\Gamma_+} &\leq \left( \int_{\Gamma_+} |\beta \cdot n|(\eta^-)^2 \right)^{\frac{1}{2}} \left( \int_{\Gamma_+} |\beta \cdot n|(\xi^-)^2 \right)^{\frac{1}{2}} \\
&\leq \left( \sum_K C_K(d, s, \beta) \frac{h_K^{s-\frac{1}{2}}}{(N+1)^{s-\frac{1}{2}}} |u|_{H^s(K)} \right) \|\|\xi\|\|
\end{aligned}
$$

Secondly, using Lemma(5.7) and applying the Cauchy-Schwarz inequality a second time leads to

$$
\begin{aligned}
\sum_K (|\beta \cdot n|\eta^-, \xi^- - \xi^+)_{\partial K_- \backslash \Gamma} &\leq \left( \sum_K \int_{\partial K_-} |\beta \cdot n|(\eta^-)^2 \right)^{\frac{1}{2}} \left( \sum_K \int_{\partial K_-} |\beta \cdot n|[\xi]^2 \right)^{\frac{1}{2}} \\
&\leq \left( \sum_K \int_{\partial K} |\beta \cdot n|(\eta^-)^2 \right)^{\frac{1}{2}} \|\|\xi\|\| \\
&\leq \left( \sum_K C_K(d, s, \beta) \frac{h_K^{s-\frac{1}{2}}}{(N+1)^{s-\frac{1}{2}}} |u|_{H^s(K)} \right) \|\|\xi\|\|
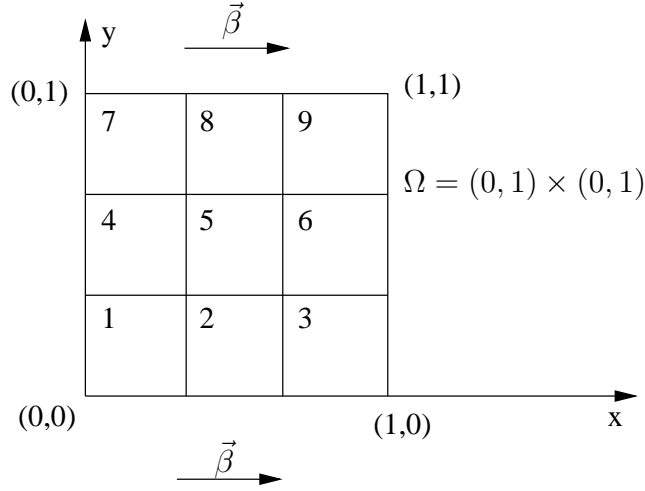\end{aligned}
$$

Figure 3: Order to compute the approximation on the elements

so that

$$\||\xi\||^2 \leq \Big( \sum_K C_K(d,s,\beta) \frac{h_K^{s-\frac{1}{2}}}{(N+1)^{s-\frac{1}{2}}} |u|_{H^s(K)} \Big) \||\xi\||$$

and

$$\||\xi\|| \leq \sum_K C_K(d,s,\beta) \frac{h_K^{s-\frac{1}{2}}}{(N+1)^{s-\frac{1}{2}}} |u|_{H^s(K)}$$

In addition, by lemma(5.6)

$$\||\eta\|| \leq C_K(d,s,\beta) \frac{h_K^{s-\frac{1}{2}}}{N^{s-\frac{1}{2}}} |u|_{H^s(K)}$$

so that

$$
\begin{aligned}
\||u - u_h\|| \;\; &\leq\;\; \||\xi\|| + \||\eta\|| \\[6pt]
&\leq\;\; \sum_K C_{1,K}(d,s,\beta) \frac{h_K^{s-\frac{1}{2}}}{(N+1)^{s-\frac{1}{2}}} |u|_{H^s(K)} + C_{2,K}(d,s,\beta) \frac{h_K^{s-\frac{1}{2}}}{N^{s-\frac{1}{2}}} |u|_{H^s(K)} \\[6pt]
&\leq\;\; \sum_K C_K(d,s,\beta) \frac{h_K^{s-\frac{1}{2}}}{N^{s-\frac{1}{2}}} |u|_{H^s(K)} \\[6pt]
&\leq\;\; C(d,s,\beta) \frac{h^{s-\frac{1}{2}}}{N^{s-\frac{1}{2}}} |u|_{H^s(\Omega)}
\end{aligned}
$$

$\square$ Theorem(5.8)

# 6 Numerical Results

In this section, the Discontinuous Galerkin method is tested numerically. The principal aim is to verify the convergence result, developed in the previous section.

## 6.1 The code

The Discontinuous Galerkin method is implemented as Matlab code. You will find it in Appendix A.

On each element, the linear system (12) is solved. As finite dimensional space $V_h$, the polynomial space $Q^N(K)$, the set of all tensor-product polynomials on $K$ of degree N in each coordinate direction is chosen. As basis of the polynomial space $Q^N(K)$, the tensor product of the Legendre's polynomials is chosen. This means that each basis $\varphi_i(x, y)$ is the product of a Legendre's polynomial in $x$ and one in $y$, each of maximum degree $N$. The Legendre's polynomials belong to Jacobi polynomials, they are orthogonal for the non weighted $L^2$-scalar product. For more details about the Legendre's polynomials, see section(8).

For computing the matrix $A$ and the vector $\vec{f}$, symoblic calculation is used. In fact, Matlab can use the Maple commands. This allows us to calculate the integrals exactly and no numerical integration is needed. For more details about Matlab using Maple, see the Matlab documentation. To solve the linear system, the GMRes-Algorithm is used.

Then, in the case of the test problems, $\|\|u - u_h\|\|$ is calculated where $u$ denotes the exact solution defined by (24) and $u_h$ is the corresponding approximation.

## 6.2 Test Problem

A two-dimensional problem is chosen. The domain is the rectangle $\Omega = (0, 1) \times (0, 1) \subset \mathbb{R}^2$ and as exact solution, the following function is taken

$$u(x, y) = \big((x - 1) + e^{-x}\big)y. \tag{24}$$

$\vec{\beta}$ is chosen as $\vec{\beta} = (1 \ 0)^T$, so that the differential equation is

$$\vec{\beta} \cdot \nabla u(x, y) + u(x, y) = xy = f(x, y)$$

$\Omega$ is divided in $N^2$ sub-rectangles (the elements), where $N$ is the number of sub-intervals of $[0, 1]$. The order to compute the approximation on the elements is illustrated in Figure(3). The boundary condition on $\Gamma_-$ is

$$u(\vec{x}) = 0 \qquad \forall \, \vec{x} \in \Gamma_-$$

It is obvious that $u \in C^\infty(\Omega)$ and $u \in H^k(\Omega) \ \forall k \geq 0$. Due to this regularity of the solution $u$, $s = \min(N + 1, k) = N + 1$ and the convergence result is in this case

$$\|\|u - u_h\|\| \leq C(d, N, \beta) \frac{h^{N+\frac{1}{2}}}{N^{N+\frac{1}{2}}} |u|_{H^{N+1}(\Omega)} \tag{25}$$

Note that the constant $C$ is also depending on N. In figure(4) we plot the logarithm of the DG-norm against the logarithm of $h$ resp. $\frac{1}{N}$ for each fixed $N$ resp. $h$. These two graphics are commented in the two following subsections.
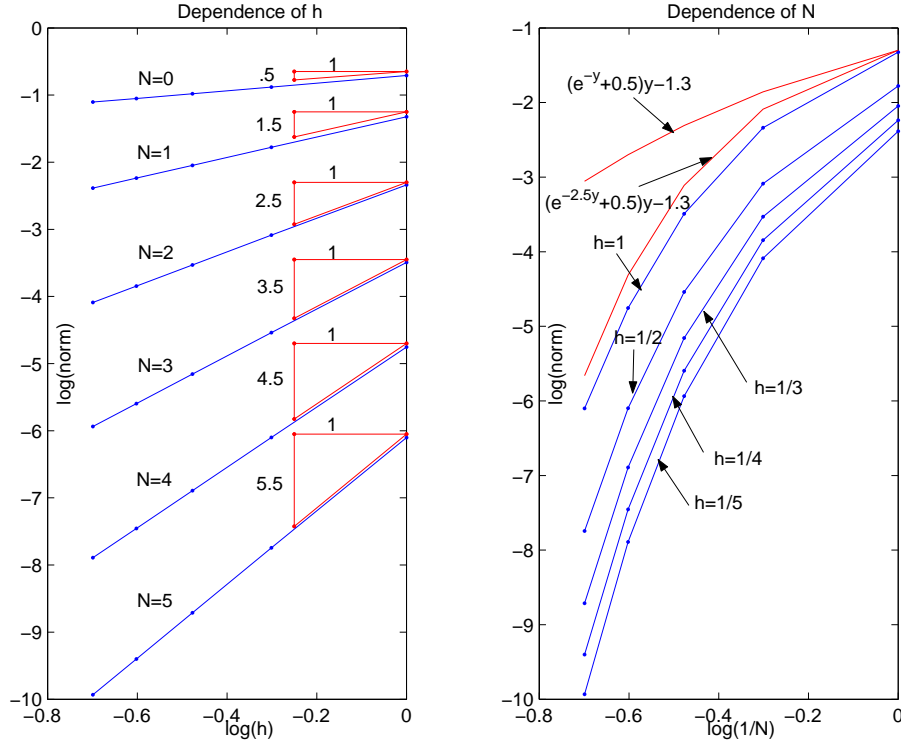
Figure 4: h- resp. N-refinement in the case of the test problem

### 6.2.1   $h$-refinement

In this case, we study the behavior of the error while varying $h$ for each fixed $N$. Fixing $N$ in (25) leads to the following convergence in $h$:

$$|||u - u_h||| \leq C_N h^{N+\frac{1}{2}}$$

In the case of a such behavior of the error, means

$$E = C_N h^{N+\frac{1}{2}} \tag{26}$$

we expect a straight line in the log-log diagram. In fact, taking the logarithm of (26) leads to

$$\log(E) = (N + \frac{1}{2}) \log(h) + \log(C_N)$$

Effecting the transformation of variables

$$y = \log(h) \Leftrightarrow h = e^{-y} \tag{27}$$

leads to

$$\log(E) = (N + \frac{1}{2})y + \log(C_N).$$

Therefor the expected straight line. As you can see in the first graphic of figure(4) we observe exactly a such behavior, i.e. the slope of the straight line varies for each $N$. So that $|||u - u_h|||$ converge at the rate $O(h^{N+\frac{1}{2}})$ as $h$ tends to zero for each fixed $N$. As you can observe, the bound is attained, means the error estimation is optimal.

### 6.2.2   $N$-refinement

In this case, $h$ is fixed. So that the convergence result gets

$$|||u - u_h||| \leq C_{h,N} N^{-(N+\frac{1}{2})} |u|_{H^{N+1}(\Omega)}$$

We neglect the $N$-dependence of the constant $C_{h,N}$ and the variation of $h^{N+\frac{1}{2}}$. Assuming the following error behavior

$$E = CN^{-(N+\frac{1}{2})} |u|_{H^{N+1}(\Omega)} \tag{28}$$

we expect the following behavior in the log-log diagram

$$\log(E) = \left(e^{-y} + \frac{1}{2}\right)y + \log(C|u|_{H^{N+1}(\Omega)}) \tag{29}$$

In fact, taking the logarithm of (28) leads to

$$\log(E) = \left(N + \frac{1}{2}\right)\log(N^{-1}) + \log(C|u|_{H^{N+1}(\Omega)})$$

Effecting the same transformation of variables as above, i.e. effecting the transformation defined by (27), leads to (29). This curve is plotted in the second graphic of figure(4).
Motivated by this calculus, we model the error by the following function.

$$f(y) = (e^{-C_1 y} + \frac{1}{2})y + C_2.$$

As you can remark, a supplementary constant $C_1$ is added and the N-dependence of $|u|_{H^{N+1}(\Omega)}$ is neglected. In the second graphic of figure(4), a such function with $C_1 = 2.5$ and $C_2 = -1.3$ is also plotted.
As predicted, we can also observe the exponential decrease of the error. Note that the exponential decrease is only due to the regularity of the solution. For example, in the case where $u \in H^1(\Omega)$, we have that $s = min(N + 1, 1) = 1, \forall N \geq 0$ and hence the convergence result becomes

$$|||u - u_h||| \leq C(d, \beta) \frac{h^{\frac{1}{2}}}{N^{\frac{1}{2}}} |u|_{H^1(\Omega)} \tag{30}$$

i.e. a convergence rate of $O\left((N^{-1})^{\frac{1}{2}}\right)$ of $|||u - u_h|||$ as $N^{-1}$ tends to zero.

## 6.3   Conclusions of the numerical tests

In both cases, the $h$-refinement and $N$-refinement, the error estimations are optimal. In the case where the solution is sufficient regular, it is more useful to refine in $N$ than in $h$, because of the exponential decrease of the error in the $N$-refinement. In the case of a weak regularity of the solution, we have no numerical results, but a behavior of the error described in (30) is expected.

# 7  Conclusion

In this section, the advantages and disadvantages of the Discontinuous Galerkin method are discussed. An advantage is that the global problem can be decoupled in a lot of small problems. For each element, a linear system has to be solved.

In addition to that, it is also easy to use spectral methods. Since the approximation does not have to be continuous, one can take an arbitrary basis of polynomials for each element. But the fact that the approximation is discontinuous implies also a larger amount of unknowns. One should compare in a further project the Continuous and Discontinuous Galerkin method in the case of the same degree of freedom.

The fact that the approximation is discontinuous can be considered as an advantage or a disadvantage. If the exact solution is continuous, then the approximation will be discontinuous, but converge to the continuous solution. On the other hand if the exact solution is discontinuous, and with a good choice of the computational mesh one may catch better the effects of the disconinuous solution than with the Continuous Galerkin method.

It should also be mentioned, that no oscillations are observed.

# 8  Introduction to Legendre's polynomials

This section is a short introduction to the family of Legendre's polynomials. The proofs of the theorems are not presented.

**Definition 8.1** *Legendre's polynomial of degree N is the polynomial defined by*

$$P_N(x) = \frac{1}{N!2^N}\frac{d^N}{dx^N}(x^2-1)^N \qquad \forall x \in [-1,1]$$

**Proposition 8.2** *$P_N$ has exactly N different real zeros all included in $(-1,1)$*

**Definition 8.3** *The zeros of the Legendre's polynomial are called the Gauss points.*

**Theorem 8.4**     *1. The Legendre's polynomials $P_1, \ldots, P_N$ build a basis of $P^N([-1,1])$*

*2. The Legendre's polynomials satisfies*

$$\int_{-1}^{1} P_n(x)P_m(x)dx = \left\{ \begin{array}{ll} 0 & if \quad n \neq m \\ \frac{2}{2n+1} & if \quad n = m \end{array} \right.$$

Note that if an orthogonal basis of $P^N([-1,1])$ is given, it is simple to construct an orthogonal basis of $P^N([a,b])$ by a simple affine transformation.

**Theorem 8.5** *If $f : [a,b] \to \mathbb{R}$ is a continuous function and let be $\{\Phi_0, \ldots, \Phi_N\}$ a set of orthonormal polynomials on $[a,b]$ where $\Phi_k$ is of degree k. Let be*

$$c_j = \int_a^b f(x)\Phi_j(x)dx$$

*Then, the polynomial defined by*

$$P = \sum_{j=0}^{N} c_j \Phi_j(x)$$

*is such that*

$$\|f - P\|_{L^2([a,b])} \leq \|f - Q\|_{L^2([a,b])} \qquad \forall Q \in P^N([a,b])$$

# 9 Introduction to numerical integration

This section is an introduction to numerical integration. Like in the previous section, only the results are presented.

**Problem** Let be $f \in C([a,b])$. One would like to approximate the functional

$$I(f) = \int_a^b f(x)\rho(x)dx \qquad \forall f \in C([a,b]) \tag{31}$$

where $\rho$ is a weight function. Let be $a \le x_0 < x_1 < \ldots < x_N \le b$ $N+1$ points in $[a,b]$.

**Definition 9.1**

$$I_{N+1}(f) = \sum_{j=0}^{N} f(x_j)\omega_j \tag{32}$$

*is called a quadrature formula where $\omega_j \in \mathbb{R}$ are called the weights. $\{x_0, \ldots, x_N\}$ are called the integration points.*

**Definition 9.2** *A quadrature formula $J$ is called exact of degree $N$ if*

- $J(f) = I(f) \qquad \forall f \in P^N$

- $J(x^{N+1}) \ne I(x^{N+1})$

**Proposition 9.3** *Let be $I$ defined by (31) and $I_{N+1}$ by (32). This formula is exact of degree $2N+1$ if*

- $\omega_i = \int_a^b L_i(x)\rho(x)$ *where $L_i$ is the i-th component of the Lagrange basis*

- *The integration points are the zeros of $P_{N+1}$, where $P_{N+1}$ belongs to the family of orthogonal polynomials on $[a,b]$.*

**Proposition 9.4** *Let be $x_0, \ldots, x_N$ $N+1$ distinct points and*

$$\omega_j = \int_a^b L_j(x)\rho(x)dx$$

*Then $I_{N+1}$ is at least exact of degree $N$ for all $f \in C^{N+1}$.*

# Acknowledgements

# References

[1] A. Ern, J.-L. Guermond: *Elément finis: théorie, applications, mise en œuvre*, Springer (2002)

[2] P. Houston, Ch. Schwab, E. Süli: *Discontinuous hp-finite element methods for advection-diffusion-reaction problems*, 2002
http://e-collection.ethbib.ethz.ch/show?type=incoll&nr=301

# 10 Appendix A

```
0001 function Norms=DG(a,b,Nx,Ny,mu,N,f,gcoeff)
0002
0003 % test arguments
0004 s=size(gcoeff);
0005 if (Ny ~= s(1))
0006     errordlg('N+1 ~= length(gcoeff)');
0007 end
0008 if (N+1 ~= s(2))
0009     errordlg('N+1 ~= length(gcoeff)');
0010 end
0011
0012 % load maple orthogonal polynomials
0013 maple('with','orthopoly')
0014
0015 % define constants
0016 Beta=[1 0]';
0017 hx=a/Nx;
0018 hy=b/Ny;
0019 Nel=Ny*Nx;
0020 if N<2
0021     numb=1;
0022 else
0023     numb=5;
0024 end
0025 h1 = hx/numb;
0026 h2 = hy/numb;
0027 counter=1;
0028
0029 % Get order of triangles and other bouandary informations
0030 % seq(1,:) = numberof element
0031 % seq(2,:) = x-coordinate
0032 % seq(3,:) = y-coordinate
0033 % seq(4,:) = number of the element which is on the left (if on the boundary=-1)
0034 % seq(5,:) = if on the bottom-boundary
0035 % seq(6,:) = if on the top-boundary
0036 % seq(7,:) = if on the right-boundary
0037 seq=zeros(5,Nel);
0038 seq(1,1:Nel) = [1:1:Nel];
0039 if(Nel>1)
0040     for i=0:Nel-1
0041         % x0
0042         seq(2,i+1) = mod(i,Nx)*hx;
0043         % y0
0044         seq(3,i+1) = (i-mod(i,Nx))/Nx*hy;
0045         % previous
0046         if (mod(i+1,Nx)==1)
0047             seq(4,i+1) = -1;
0048         else
0049             seq(4,i+1) = i;
0050         end
0051         if i<Nx
0052             seq(5,i+1) = 1;
0053         end
0054         if i>=(Nel-Nx)
0055             seq(6,i+1) = 1;
0056         end
0057         if mod(i+1,Nx)==0
0058             seq(7,i+1) = 1;
0059         end
0060     end
0061 else
0062     seq(4,1) = -1;
0063     seq(5,1) = 1;
0064     seq(6,1) = 1;
0065     seq(7,1) = 1;
0066 end
0067
0068 % define matrices
0069 U = zeros((N+1)^2,Nel);
0070 A = zeros((N+1)^2,(N+1)^2);
0071 Ak = zeros((N+1)^2,(N+1)^2);
0072 Bk = zeros((N+1)^2,(N+1)^2);
```

```
0073 Ck = zeros((N+1)^2,(N+1)^2);
0074 Aint = zeros(N+1,N+1);
0075 Ufunc = cell(Nel,1);
0076
0077 for i=0:N
0078     Aint(i+1,i+1) = 2/(2*i+1);
0079 end
0080
0081 % for the norm
0082 if (strcmp(f,'x*y') | strcmp(f,'(10*cos(10*x)+sin(10*x))*y^3') & mu==1)
0083     if (strcmp(f,'x*y'))
0084         exactsol = '((x-1)+exp(-x))*y';
0085     else
0086         exactsol = '(sin(10*x))*y^3';
0087     end
0088     Lnorm2=0;
0089     DGnorm2=0;
0090 else
0091     'error norm calculation is not possible, exact solution is not known'
0092     %warndlg('error norm calculation is not possible, exact solution is not known');
0093 end
0094
0095 % Construction of Ak
0096 Tk = [[hx/2 0];[0 hy/2]];
0097 for i=0:N
0098     for j=0:N
0099         for k=0:N
0100             for l=0:N
0101                 Iind = i*(N+1)+j+1;
0102                 Jind = k*(N+1)+l+1;
0103                 % Ak
0104                 Ak(Iind,Jind) = mu*det(Tk)*Aint(i+1,k+1)*Aint(j+1,l+1);
0105                 % Bk
0106                 Bk(Iind,Jind) = det(Tk)*dot(Beta,inv(Tk)'*...
0107                     [Aint(j+1,l+1)*getInt(i,k); Aint(i+1,k+1)*getInt(j,l)]);
0108                 % Ck
0109                 val1 = str2num(maple('P',i,-1));
0110                 val2 = str2num(maple('P',k,-1));
0111                 Ck(Iind,Jind) = hy/2*val1*val2*Aint(j+1,l+1);
0112             end
0113         end
0114     end
0115 end
0116 Ak = Ak + Bk + Ck;
0117
0118 % loop on the elements
0119 for k=1:Nel
0120     % get f
0121     bk = [seq(2,k)+hx/2 seq(3,k)+hy/2]';
0122     x0 = seq(2,k);
0123     x0str = maple('eval',x0);
0124     x0strPlus = maple('eval',x0+hx);
0125     y0 = seq(3,k);
0126     y0str = maple('eval',y0);
0127     y0strPlus = maple('eval',y0+hy);
0128     invTransX = ['2/(' maple('eval',hx) ')*(x-' maple('eval',x0) ')-1'];
0129     invTransY = ['2/(' maple('eval',hy) ')*(y-' maple('eval',y0) ')-1'];
0130     if (seq(4,k)==-1)
0131         coeff = gcoeff(counter,:);
0132     else
0133         coeff = outflow(:,seq(4,k));
0134     end
0135     for i1=0:N
0136         % First Term
0137         pi1 = maple('P',i1,invTransX);
0138         func = ['(' f ')*(' pi1 ')'];
0139         str = maple('int',func,['x=' x0str '..' x0strPlus]);
0140         % Second Term
0141         pi1atmin1 = str2num(maple('P',i1,-1));
0142         for i2=0:N
0143             % First Term
0144             pi2 = maple('P',i2,invTransY);
0145             func = ['(' pi2 ')*(' str ')'];
0146             str2 = maple('int',func,['y=' y0str '..' y0strPlus]);
```

```
0147                % Second Term
0148                second = abs(Beta'*[-1 0]')*hy/2*pi1atmin1*coeff(i2+1)*2/(2*i2+1);
0149                % Sum
0150                fk(i1*(N+1)+i2+1) = str2num(str2) + second;
0151            end
0152        end
0153
0154        if (seq(4,k)==-1)
0155            counter=counter+1;
0156        end
0157
0158        % get solution
0159        U(:,k) = gmres(Ak,fk',20,10e-15,500);%Ak\fk';
0160
0161        % Get grid for evaluation
0162        x = [x0:h1:x0+hx];
0163        y = [y0:h2:y0+hy];
0164        [X(:,:,k) Y(:,:,k)] = meshgrid(x,y);
0165
0166        % get outflow
0167        outflow(:,k) = zeros(N+1,1);
0168
0169        % get solution function in form of string
0170        % loop over x-basis
0171        for i=1:N+1
0172            % loop over y-basis
0173            for j=1:N+1
0174                % get func
0175                if (i==1 & j==1)
0176                    Ufunc(k) = {['(' num2str(U(1,k),15) '*(' maple('P',0,invTransX)...
0177                                 ')*(' maple('P',0,invTransY) '))']};
0178                else
0179                    Ufunc(k) = {[char(Ufunc(k)) '+(' num2str(U((i-1)*(N+1)+j,k),15)...
0180                                 '*(' maple('P',i-1,invTransX) ')*(' maple('P',j-1,invTransY) '))']};
0181                end
0182                % get i-th outflow-coeff for element k
0183                outflow(i,k) = outflow(i,k) + str2num(maple('P',j-1,1))*U((j-1)*(N+1)+i,k);
0184            end
0185        end
0186
0187        % evaluate solution function on grid points for visualization
0188        s = size(X(:,:,k));
0189        for i=1:s(1)
0190            for j=1:s(2)
0191                x=X(i,j,k);
0192                y=Y(i,j,k);
0193                v(i,j,k) = eval(char(Ufunc(k)));
0194            end
0195        end
0196
0197        % calculates norms of this element
0198        if (strcmp(f,'x*y') | strcmp(f,'(10*cos(10*x)+sin(10*x))*y^3') & mu==1)
0199            if (strcmp(f,'x*y'))
0200                Lnorm2 = Lnorm2 + L2norm(x0,y0,hx,hy,char(Ufunc(k)),exactsol)^2;
0201                DGnorm2 = DGnorm2 + pureDGnorm(x0,y0,hx,hy,Ufunc,seq(:,k),k,exactsol)^2;
0202            else
0203                Lnorm2 = Lnorm2 + L2norm(x0,y0,hx,hy,char(Ufunc(k)),exactsol)^2;
0204                DGnorm2 = DGnorm2 + pureDGnorm(x0,y0,hx,hy,Ufunc,seq(:,k),k,exactsol)^2;
0205            end
0206        end
0207 end
0208
0209 % norm calculation
0210 if (strcmp(f,'x*y') | strcmp(f,'(10*cos(10*x)+sin(10*x))*y^3') & mu==1)
0211    Norms(1) = sqrt(Lnorm2);
0212    Norms(2) = sqrt(DGnorm2);
0213    Norms(3) = sqrt(DGnorm2+Lnorm2);
0214 else
0215    Norms = [-1 -1 -1];
0216 end
0217
0218 % visualisation
0219 figure(33)
0220 for k=1:Nel
```

```
0221      patch(surf2patch(X(:,:,k)',Y(:,:,k)',v(:,:,k)',v(:,:,k)'));
0222      shading faceted;
0223      view(3)
0224 end
0225 hold off;
0226
```