# Simple finite element methods in Python

Roland Becker

April 22, 2022

## Contents

# 1 Elliptic equation

Let $\Omega \subset \mathbb{R}^d$, $d = 2, 3$ be the computational domain. We suppose to have a disjoined partition of its boundary: $\partial\Omega = \Gamma_D \cup \Gamma_N \cup \Gamma_R$. We consider the second-order elliptic equation for the scalar unknown $u$ with coefficients

$$A(x) \in \mathbb{R}^{d \times d}, \quad b(x) \in \mathbb{R}^d, \quad c(x) \in \mathbb{R}, \quad \alpha \in \mathbb{R}.$$

---

**Elliptic Equation (strong formulation)**

$$\begin{cases} -\operatorname{div}(A\nabla u) + b \cdot \nabla u + cu = f & \text{in } \Omega \\ u = u^D & \text{on } \Gamma_D \\ \left|b_n^-\right| u + n^{\mathsf{T}} A\nabla u = q^N & \text{on } \Gamma_N \\ \alpha u + \left|b_n^-\right| u + n^{\mathsf{T}} A\nabla u = u^R & \text{on } \Gamma_R \end{cases} \tag{1.1}$$

For the standard Neumann condition, we need $b_n^- = 0$, i.e. $\Gamma_N$ is part of the outflow boundary.

---

## 1.1 Standard weak formulation

---

**Elliptic Equation (weak formulation)**

Let $H_\phi^1 := \left\{ u \in H^1(\Omega) \mid u\big|_{\Gamma_D} = \phi \right\}$. The primal weak formulation looks for $u \in H_{u^D}^1$ such that for all $\phi \in H_0^1(\Omega)$

$$\int_\Omega A\nabla u \cdot \nabla \phi + \int_\Omega (b \cdot \nabla u)\phi + \int_\Omega cu\phi + \int_{\Gamma_R} \alpha u\phi + \int_{\Gamma_R \cup \Gamma_N} \left|b_n^-\right| u\phi = \int_\Omega f\phi + \int_{\Gamma_N} q^N\phi + \int_{\Gamma_R} u^R\phi \tag{1.2}$$

---

We can derive (1.2) from (1.1) by the divergence theorem

$$\int_\Omega \operatorname{div} \vec{F} = \int_{\partial\Omega} \vec{F}_n \quad \overset{F \to F\phi}{\Longrightarrow} \quad \int_\Omega (\operatorname{div} \vec{F})\phi = -\int_\Omega \vec{F} \cdot \nabla\phi + \int_{\partial\Omega} \vec{F}_n\phi,$$

which gives with $\vec{F} = A\nabla u$ and $\phi \in H_0^1(\Omega)$

$$\int_\Omega \operatorname{div}(A\nabla u)\,\phi = -\int_\Omega A\nabla u \cdot \nabla\phi + \int_{\Gamma_N \cup \Gamma_R} n^{\mathsf{T}} A\nabla u\phi.$$

It follows that for a sufficiently smooth solution $u$ of (1.2) we have

$$\int_\Omega \left(-\operatorname{div}(A\nabla u) + b \cdot \nabla u + cu - f\right)\phi = \int_{\Gamma_N} \left(q^N - n^{\mathsf{T}} A\nabla u - \left|b_n^-\right| u\right)\phi + \int_{\Gamma_R} \left(u^R - n^{\mathsf{T}} A\nabla u - \alpha u - \left|b_n^-\right| u\right)\phi$$

Taking first $\phi$ vanishing on the whole boundary, we find that (1.1) is satisfied almost everywhere in $\Omega$. Second, we recover the Neumann and Robin boundary conditions. Notice that the Dirichlet condition is imposed by the space.

Denoting the bilinear form on the left of (1.2) by $a$, we have in case $u^D = 0$, taking $\phi = u$, and with

$$\int_\Omega (b \cdot \nabla u)\phi + \int_\Omega u(b \cdot \nabla \phi) + \int_\Omega (\operatorname{div} b)u\phi = \int_{\partial\Omega} b_n u\phi \quad \Rightarrow \quad \int_\Omega (b \cdot \nabla u)\phi = -\int_\Omega \frac{\operatorname{div} b}{2}u^2 + \int_{\partial\Omega} \frac{b_n}{2}u^2$$

$$a(u, u) = \int_\Omega A\nabla u \cdot \nabla u + \int_\Omega (c - \frac{\operatorname{div} b}{2})u^2 + \int_{\Gamma_R \cup \Gamma_N} \frac{|b_n|}{2}u^2 + \int_{\Gamma_R} \alpha u^2.$$

We obtain coercivity under the standard assumptions

$$\xi^T A(x)\xi \geqslant \alpha_0 |\xi|^2, \quad c - \frac{\operatorname{div} b}{2} \geqslant 0, \quad \alpha \geqslant 0. \tag{1.3}$$

## 1.2  Mixed weak formulation

Introducing the flux as an unknown

$$q := A\nabla u, \quad q \in Q_\phi := \left\{ q \in H(\operatorname{div}, \Omega) \,\middle|\, q|_{\Gamma_N} = \phi \right\}, \quad C := A^{-1}$$

$$\begin{cases} (u, q) \in L^2(\Omega) \times Q_{q^N} \\[2mm] \displaystyle\int_\Omega Cq \cdot p + \frac{1}{\alpha}\int_{\Gamma_R} q_n p_n + \int_\Omega u \operatorname{div} p \;=\; \int_{\Gamma_D} u^D p_n + \frac{1}{\alpha}\int_{\Gamma_R} u^R p_n \qquad \forall p \in Q_0 \\[2mm] \displaystyle\int_\Omega v \operatorname{div} q - \int_\Omega (b \cdot \nabla u + cu)v \qquad\qquad = -\int_\Omega fv \qquad\qquad\qquad \forall v \in L^2(\Omega). \end{cases} \tag{1.4}$$

Integration by parts in the first equation gives

$$\int_\Omega (Cq - \nabla u) \cdot p = \int_{\Gamma_D} (u^D - u)p_n + \int_{\Gamma_R} (\frac{1}{\alpha}u^R - u - \frac{1}{\alpha}q_n)p_n$$

## 1.3  Boundary conditions

### 1.3.1  Nitsche's method

$$\begin{cases} u_h \in V_h: \quad a_\Omega(u_h, \phi) + a_{\partial\Omega}(u_h, \phi) = l_\Omega(\phi) + l_{\partial\Omega}(\phi) \quad \forall \phi \in V_h \\[2mm] a_\Omega(v, \phi) := \displaystyle\int_\Omega \mu\nabla u \cdot \nabla \phi \\[2mm] a_{\partial\Omega}(v, \phi) := \displaystyle\int_{\Gamma_D} \frac{\gamma\mu}{h}u\phi - \int_{\Gamma_D} \mu\left(\frac{\partial u}{\partial n}\phi + u\frac{\partial \phi}{\partial n}\right) \\[2mm] l_\Omega(\phi) := \displaystyle\int_\Omega f\phi, \quad l_{\partial\Omega}(\phi) = \int_{\Gamma_D} \mu u^D\left(\frac{\gamma}{h}\phi - \frac{\partial \phi}{\partial n}\right) \end{cases} \tag{1.5}$$

Let $-\operatorname{div}(\mu\nabla z) = 0$ and $z|_{\Gamma_D} = 1$ and $z|_{\Gamma_N} = 0$. Then

$$\int_\Omega \mu\nabla u \cdot \nabla z - \int_\Omega fz = \int_\Omega (\mu\nabla u \cdot \nabla z + \operatorname{div}(\mu\nabla u)z) = \int_{\Gamma_D} \mu\frac{\partial u}{\partial n}.$$

4

Now, if $z_h \in V_h$ such that $z - z_h \in H_0^1(\Omega)$

$$\int_\Omega \mu \nabla(u - u_h) \cdot \nabla(z - z_h) = \int_\Omega f(z - z_h) - \int_\Omega \mu \nabla u_h \cdot \nabla(z - z_h)$$

$$= \int_\Omega fz - \int_\Omega \mu \nabla u_h \cdot \nabla z + \int_\Omega \mu \nabla u_h \cdot \nabla(z - z_h) - \int_\Omega f z_h$$

$$= -\int_{\Gamma_D} \mu \frac{\partial u}{\partial n} + \int_\Omega \mu \nabla(u - u_h) \cdot \nabla z + \int_{\Gamma_D} \mu(u^D - u_h) \left( \frac{\gamma}{h} z_h - \frac{\partial z_h}{\partial n} \right) + \int_{\Gamma_D} \mu \frac{\partial u_h}{\partial n}$$

$$= \int_{\Gamma_D} \mu \frac{\partial u_h}{\partial n} + \int_{\Gamma_D} (u^D - u_h) \frac{\mu \gamma}{h} - \int_{\Gamma_D} \mu \frac{\partial u}{\partial n} + \int_{\Gamma_D} \mu(u - u_h) \frac{\partial(z - z_h)}{\partial n},$$

so we get a possibly second-order approximation of the flux by

$$F_h := \int_{\Gamma_D} \mu \frac{\partial u_h}{\partial n} + \int_{\Gamma_D} (u^D - u_h) \frac{\mu \gamma}{h}. \tag{1.6}$$

## 1.4 Computation of the matrices for $\mathcal{P}_h^1(\Omega)$

For the convection, we suppose that $\vec{v} \in \mathcal{RT}_h^0(\Omega)$ and let for given $K \in \mathcal{K}_h$ $\vec{v} = \sum_{k=1}^{d+1} v_k \Phi_k$. Using

$$x_k = x_{S_k}^K, \quad h_k = h_{S_k}^K, \quad \sigma_k = \sigma_{S_k}^K, n_k = n_{S_k}$$

we compute

$$\int_K \lambda_j \vec{v} \cdot \nabla \lambda_i = \sum_{k=1}^{d+1} v_k \int_K \lambda_j \Phi_k \cdot \nabla \lambda_i$$

$$\int_K \lambda_j \Phi_k \cdot \nabla \lambda_i = -\frac{\sigma_k \sigma_i}{h_k h_i} \int_K \lambda_j (x - x_k) \cdot n_i = -\frac{\sigma_k \sigma_i}{h_k h_i} \sum_{l=1}^{d+1} (x_l - x_k) \cdot n_i \int_K \lambda_j \lambda_l$$

## 2 Heat equation

Let $\Omega \subset \mathbb{R}^d$, $d = 2, 3$ be the computational domain. We suppose to have a disjoined partition of its boundary: $\partial\Omega = \Gamma_D \cup \Gamma_N \cup \Gamma_R$. We consider the parabolic equation for the temperature $T$, heat flux $\vec{q}$ and heat release $\dot{q}$

**Heat equation (strong formulation)**

$$
\begin{cases}
\vec{q} = -k\nabla T \\[2mm]
\rho C_p \dfrac{dT}{dt} + \operatorname{div}(\vec{v}T) + \operatorname{div}\vec{q} = \dot{q} & \text{on } \Omega \\[2mm]
T = T^D & \text{in } \Gamma_D \\[2mm]
k\dfrac{\partial T}{\partial n} = q^N & \text{on } \Gamma_N \\[2mm]
c_R T + k\dfrac{\partial T}{\partial n} = q^R & \text{on } \Gamma_R
\end{cases}
\tag{2.1}
$$

**Heat equation (weak formulation)**

Let $H_\phi^1 := \left\{ T \in H^1(\Omega) \,\middle|\, T|_{\Gamma_D} = \phi \right\}$. The standard weak formulation looks for $T \in H_{T^D}^1$ such that for all $\phi \in H_0^1(\Omega)$

$$
\int_\Omega \rho C_p \frac{dT}{dt}\phi - \int_\Omega \vec{v}T \cdot \nabla\phi + \int_\Omega k\nabla T \cdot \nabla\phi + \int_{\Gamma_R} c_R T\phi + \int_{\Gamma_R \cup \Gamma_N} \vec{v}_n T\phi = \int_\Omega \dot{q}\phi + \int_{\Gamma_R} q^R \phi \tag{2.2}
$$

We can derive (2.2) from (2.1) by the divergence theorem

$$
\int_\Omega \operatorname{div}\vec{F} = \int_{\partial\Omega} \vec{F}_n \quad \overset{F \to F\phi}{\Longleftrightarrow} \quad \int_\Omega (\operatorname{div}\vec{F})\phi = -\int_\Omega \vec{F} \cdot \nabla\phi + \int_{\partial\Omega} \vec{F}_n \phi,
$$

which gives with $\vec{F} = \vec{v} + \vec{q}$

$$
\int_\Omega \operatorname{div}(\vec{v} + \vec{q})\,\phi = -\int_\Omega \vec{v} \cdot \nabla\phi + \int_\Omega k\nabla T \cdot \nabla\phi + \int_{\partial\Omega} \vec{F}_n \phi.
$$

Using that $\phi$ vanishes on $\Gamma_D$ we have

$$
\int_{\partial\Omega} \vec{F}_n \phi = \int_{\Gamma_N \cup \Gamma_R} \vec{F}_n \phi = \int_{\Gamma_N \cup \Gamma_R} \vec{v}_n \phi + \int_{\Gamma_N \cup \Gamma_R} \vec{q}_n \phi,
$$

and then with the different boundary conditions, we find

$$
\int_{\Gamma_N \cup \Gamma_R} \vec{q}_n \phi = \int_{\Gamma_D} q^N \phi + \int_{\Gamma_R} \left( q^R - c_R T \right)\phi
$$

# 3 Stokes problem

Let $\Omega \subset \mathbb{R}^d$, $d = 2, 3$ be the computational domain. We suppose to have a disjoined partition of its boundary: $\partial\Omega = \Gamma_D \cup \Gamma_N \cup \Gamma_R \cup \Gamma_P$, the outward unit normal is denotes by $n$. We denote by $\pi_{n\perp} := (I - nn^\mathsf{T}) \in \mathbb{R}^{(d-1)\times d}$ the projection on the tangent plane.

$$
\begin{cases}
\begin{aligned}
-\operatorname{div}(\mu\nabla v) + \nabla p &= f \quad \text{in } \Omega \\
\operatorname{div} v &= g \quad \text{in } \Omega \\
v &= v^D \quad \text{in } \Gamma_D,
\end{aligned} \\[2mm]
\left.\begin{aligned}
\pi_{n\perp}\mu\frac{\partial v}{\partial n} &= 0 \\
\mu\frac{\partial v}{\partial n}\cdot n - p &= -p^N
\end{aligned}\right\} \quad \text{in } \Gamma_N, \\[4mm]
\left.\begin{aligned}
v_n &= v_n^R \\
\pi_{n\perp}\left(\lambda_R v + \mu\frac{\partial v}{\partial n}\right) &= \pi_{n\perp}g^R
\end{aligned}\right\} \quad \text{in } \Gamma_R, \\[4mm]
\left.\begin{aligned}
\pi_{n\perp}v &= \pi_{n\perp}v^P \\
\left(\mu\frac{\partial v}{\partial n}\cdot n - p\right) &= -p^P
\end{aligned}\right\} \quad \text{in } \Gamma_P.
\end{cases}
\tag{3.1}
$$

We can express the equations by means of the Cauchy stress tensor

$$
\sigma := 2\mu D(v) + \lambda\operatorname{div}(v)I - pI, \quad D(v) = \frac{1}{2}\left(\nabla v + \nabla v^\mathsf{T}\right). \tag{3.2}
$$

Due to $\operatorname{div} v = 0$, the bulk viscosity $\lambda$ is neglected.

Then the momentum balance is (with the row-wise divergence operator)

$$
-\operatorname{div}\sigma = f \quad \text{in } \Omega.
$$

The weak formulation (3.4) is based on the non-symmetric

$$
\widetilde{\sigma} := \mu\nabla v - pI, \tag{3.3}
$$

which is equivalent to using $\sigma$ for volume integrals since $A : B = \frac{A+A^\mathsf{T}}{2} : B$ for all symmetric $B \in \mathbb{R}^{d\times d}$ and any $A \in \mathbb{R}^{d\times d}$.

Using $\sigma$ in a weak formulation will in general generate different boundary conditions.

## 3.1 Weak formulation

The standard weak formulation reads

$$
\begin{cases}
V_{v^D, v_n^R, v^P} := \left\{ v \in H^1(\Omega, \mathbb{R}^d) \ \middle|\ v\big|_{\Gamma_D} = v^D \ \& \ v_n\big|_{\Gamma_R} \ \& \ \pi_{n\perp}v\big|_{\Gamma_P} = \pi v^P \right\} \\[2mm]
Q := L^2(\Omega) \quad (Q := L^2(\Omega)/\mathbb{R} \quad \text{if } |\Gamma_N \cup \Gamma_P| = 0) \\[2mm]
(v, p) \in V_{v^D, v_n^R, v^P} \times Q: \quad a_\Omega(v, p; \phi\xi) = l_\Omega(\phi, \xi) \quad \forall(\phi, \xi) \in V_{0,0,0} \times Q \\[2mm]
a_\Omega(v, p; \phi, \xi) := \displaystyle\int_\Omega \mu\nabla v : \nabla\phi - \int_\Omega p\operatorname{div}\phi + \int_\Omega \operatorname{div} v\xi + \lambda_R\int_{\Gamma_R}(v\cdot\phi - v_n\phi_n), \\[4mm]
l_\Omega(\phi, \xi) := \displaystyle\int_\Omega f\cdot\phi + \int_\Omega g\xi + \int_{\Gamma_R}(g^R\cdot\phi - g_n^R\phi_n) - \int_{\Gamma_N} p^N\phi_n - \int_{\Gamma_P} p^P\phi_n.
\end{cases}
\tag{3.4}
$$

**Lemma 3.1.** *A regular solution of the formulation (3.4) satisfies (3.1).*

*Proof.* By integration by parts we have, together with $v \cdot \phi - v_n \phi_n = \pi_{n\perp} v \cdot \phi$

$$a_\Omega(v, p; \phi, \xi) = \int_\Omega (-\mu\Delta v + \nabla p) \cdot \phi + \int_{\partial\Omega} \mu\frac{\partial v}{\partial n} \cdot \phi - \int_{\partial\Omega} p\phi_n + \int_\Omega \operatorname{div} v\xi + \lambda_R \int_{\Gamma_R} \pi_{n\perp} v \cdot \phi$$

Then the (regular) weak solution satisfies for all $\phi$

$$\int_\Omega (-\mu\Delta v + \nabla p - f) \cdot \phi = \int_{\partial\Omega} p\phi_n - \int_{\partial\Omega} \mu\frac{\partial v}{\partial n} \cdot \phi - \int_{\Gamma_N} p^N \phi_n - \int_{\Gamma_P} p^P \phi_n + \int_{\Gamma_R} \pi_{n\perp}(g^R - \lambda_R v) \cdot \phi,$$

and for all $\xi$

$$\int_\Omega (\operatorname{div} v - g)\xi = 0.$$

Taking $\phi \in H_0^1(\Omega, \mathbb{R}^d)$, the right-hand side vanishes and the density of this space in $L^2(\Omega)$ gives us

$$-\mu\Delta v + \nabla p = f, \quad \operatorname{div} v = g \quad \text{a.e. in } \Omega.$$

But this means that for general $\phi \in V_{0,0,0}$

$$\int_{\partial\Omega} p\phi_n - \int_{\partial\Omega} \mu\frac{\partial v}{\partial n} \cdot \phi - \int_{\Gamma_N} p^N \phi_n - \int_{\Gamma_P} p^P \phi_n + \int_{\Gamma_R} \pi_{n\perp}(g^R - \lambda_R v) \cdot \phi = 0$$

Decomposing the test function as

$$\phi = \phi_n n + \pi_{n\perp}\phi$$

and using the definition of $V_{0,0,0}$ we find

$$\int_{\Gamma_N} \left((p - p^N)n - \mu\frac{\partial v}{\partial n}\right) \cdot \phi + \int_{\Gamma_P} (p - p^P - \mu\frac{\partial v}{\partial n} \cdot n)\phi_n + \int_{\Gamma_R} \pi_{n\perp}(g^R - \lambda_R v - \mu\frac{\partial v}{\partial n}) \cdot \phi = 0$$

$\square$

**Proposition 3.2.** *If we use the weak formulation based on the stress tensor*

$$a_\Omega(v, p; \phi, \xi) := \int_\Omega \sigma : \nabla\phi + \int_\Omega \operatorname{div} v\xi + \lambda_R \int_{\Gamma_R} \pi_{n\perp} v \cdot \pi_{n\perp}\phi, \tag{3.5}$$

*the resulting boundary conditions are*

$$\begin{cases} v = v^D & in\ \Gamma_D, \\ \sigma n = -p^N n & in\ \Gamma_N, \\ \left. \begin{cases} \pi_{n\perp} v = \pi_{n\perp} v^P \\ (I - \pi_{n\perp})\sigma n = -p^P \end{cases} \right\} & in\ \Gamma_P, \\ \left. \begin{cases} v_n = v_n^R \\ \pi_{n\perp}\left(\lambda_R v + \sigma n\right) = \pi_{n\perp} g^R \end{cases} \right\} & in\ \Gamma_R. \end{cases}$$

$\Leftrightarrow$

$$\begin{cases} v = v^D & in\ \Gamma_D, \\ \sigma n = -p^N n & in\ \Gamma_N, \\ \left. \begin{cases} \pi_{n\perp}\mu\dfrac{\partial v}{\partial n} - \mu\omega n = 0 \\ \mu\dfrac{\partial v}{\partial n}\cdot n - p = -p^N \end{cases} \right\} & in\ \Gamma_N, \\ \left. \begin{cases} \pi_{n\perp} v = \pi_{n\perp} v^P \\ \mu\dfrac{\partial v}{\partial n}\cdot n - p = -p^P \end{cases} \right\} & in\ \Gamma_P, \\ \left. \begin{cases} v_n = v_n^R \\ \pi_{n\perp}\left(\lambda_R v + \mu\dfrac{\partial v}{\partial n} - \mu\omega n\right) = \pi_{n\perp} g^R \end{cases} \right\} & in\ \Gamma_R. \end{cases}$$

(3.6)

*Proof.* Using now

$$\int_\Omega \sigma : \nabla\phi = -\int_\Omega \operatorname{div}\sigma\cdot\phi + \int_{\partial\Omega}\sigma n\cdot\phi \tag{3.7}$$

we get in similar way as before

$$\int_{\Gamma_N}\left(p^N n + \sigma n\right)\cdot\phi + \int_{\Gamma_P}(I - \pi_{n\perp})\left(p^P n + \sigma n\right)\cdot\phi + \int_{\Gamma_R}\pi_{n\perp}\left(\lambda_R v + \sigma n - g^R\right)\cdot\phi = 0$$

We have

$$n^\mathsf{T}(\nabla v)^\mathsf{T} n = (\nabla v n)^\mathsf{T} n = \frac{\partial v}{\partial n}^\mathsf{T} n = n^\mathsf{T}\frac{\partial v}{\partial n} = n^\mathsf{T}(\nabla v)n$$

and therefore

$$n^\mathsf{T}\sigma n = n^\mathsf{T}\tilde\sigma n \quad\Rightarrow\quad (I - \pi_{n\perp})\sigma n = (I - \pi_{n\perp})\tilde\sigma n,$$
$$\pi_{n\perp}\sigma n = (I - nn^\mathsf{T})\sigma n = \sigma n - nn^\mathsf{T}\tilde\sigma n = \pi_{n\perp}\tilde\sigma n + (\sigma - \tilde\sigma)n = \pi_{n\perp}\tilde\sigma n - \mu\omega n$$

with $\omega := \frac{1}{2}(\nabla v - \nabla v^\mathsf{T})$. $\qquad\square$

## 3.2 Discretization

We use finite element spaces $V_h$ for the velocity and $Q_h$ for the pressure. One main difficulty is to obtain a stable approximation of the pressure gradient, which requires the inf-sup condition

$$\inf_{p\in Q_h\backslash\{0\}}\sup_{v\in V_h\backslash\{0\}}\frac{\int_\Omega p\operatorname{div}v}{\|v\|_V\|p\|_Q}\geqslant\gamma > 0. \tag{3.8}$$

To this end, we use the classical spaces $V_h = \mathcal{CR}_h^1(\Omega, \mathbb{R}^d)$ and $Q_h = \mathcal{D}_h^0$.

### 3.3 Implementations of Boundary condition

#### 3.3.1 Strong implementation of Dirichlet condition

We write the discrete velocity space $V_h$ as a direct sum $V_h = V_h^{int} \oplus V_h^{dir}$, with $V_h^{dir}$ corresponding to the discrete functions not vanishing on $\Gamma_D$. Splitting the matrix and right-hand side vector correspondingly, and letting $u_h^D \in V_h^{dir}$ be an approximation of the Dirichlet data $v^D$ we have the traditional way to implement Dirichlet boundary conditions:

$$\begin{bmatrix} A^{int} & 0 & -B^{int\,T} \\ 0 & I & 0 \\ B^{int} & 0 & 0 \end{bmatrix} \begin{bmatrix} v_h^{int} \\ v_h^{dir} \\ p_h \end{bmatrix} = \begin{bmatrix} f^{int} - A^{int,dir}v_h^D \\ v_h^D \\ g - B^{dir}v_h^D \end{bmatrix}. \tag{3.9}$$

As for the Poisson problem, we obtain an alternative formulation

$$\begin{bmatrix} A^{int} & 0 & -B^{int\,T} \\ 0 & A^{dir} & 0 \\ B^{int} & 0 & 0 \end{bmatrix} \begin{bmatrix} v_h^{int} \\ v_h^{dir} \\ p_h \end{bmatrix} = \begin{bmatrix} f^{int} - A^{int,dir}v_h^D \\ A^{dir}v_h^D \\ g - B^{dir}v_h^D \end{bmatrix}. \tag{3.10}$$

#### 3.3.2 Weak implementation (Nitsche's method)

Instead of modifying the discrete velocity space, we add additional terms to the bilinear and linear forms.

$$\begin{cases} (v,p) \in V_h \times Q_h: \quad a_\Omega(v,p;\phi\xi) + a_{\partial\Omega}(v,p;\phi,\xi) = l_\Omega(\phi,\xi) + l_{\partial\Omega}(\phi,\xi) \quad \forall(\phi,\xi) \in V_h \times Q_h \\[6pt] a_{\partial\Omega}(v,p;\phi,\xi) := \displaystyle\int_{\Gamma_D \cup \Gamma_P} \frac{\gamma\mu}{h} v \cdot \phi - \int_{\Gamma_D \cup \Gamma_P} \mu \left( \frac{\partial v}{\partial n} \cdot \phi + v \cdot \frac{\partial \phi}{\partial n} \right) \\[10pt] \qquad + \displaystyle\int_{\Gamma_R} \frac{\gamma\mu}{h} v_n \phi_n - \int_{\Gamma_R} \mu \left( \frac{\partial v}{\partial n} \cdot n\phi_n + v_n \frac{\partial \phi}{\partial n} \cdot n \right) \\[10pt] \qquad - \displaystyle\int_{\Gamma_P} \frac{\gamma\mu}{h} v_n \phi_n + \int_{\Gamma_P} \mu \left( \frac{\partial v}{\partial n} \cdot n\phi_n + v_n \frac{\partial \phi}{\partial n} \cdot n \right) \\[10pt] \qquad + \displaystyle\int_{\Gamma_D \cup \Gamma_R} (p\phi_n - v_n\xi) \\[10pt] l_{\partial\Omega}(\phi,\xi) = \displaystyle\int_{\Gamma_D} \mu v^D \cdot \left( \frac{\gamma}{h}\phi - \frac{\partial\phi}{\partial n} \right) - \int_{\Gamma_D} v_n^D \xi + \int_{\Gamma_R} \mu v_n^R \cdot \left( \frac{\gamma}{h}\phi_n - \frac{\partial\phi}{\partial n} \cdot n \right) - \int_{\Gamma_R} v_n^R \xi \\[10pt] \qquad + \displaystyle\int_{\Gamma_P} \mu \pi_{n^\perp} v^P \cdot \pi_{n^\perp} \left( \frac{\gamma}{h}\phi - \frac{\partial\phi}{\partial n} \right). \end{cases} \tag{3.11}$$

**Lemma 3.3.** *A regular continuous solution of the formulation (3.4) satisfies (3.11).*

*Proof.* We have already seen that a regular continuous solution satisfies for $(\phi,\xi) \in V_h \times Q_h$

$$a_\Omega(v,p;\phi,\xi) - l_\Omega(\phi,\xi) = \int_{\Gamma_D} \left( \mu\frac{\partial v}{\partial n} - pn \right) \cdot \phi + \int_{\Gamma_R} \left( \mu\frac{\partial v}{\partial n} \cdot n - p \right) \phi_n + \int_{\Gamma_P} \pi_{n^\perp}\mu\frac{\partial v}{\partial n} \cdot \phi$$

Thanks to the boundary conditions we also have

$$\int_{\Gamma_D} \mu(v^D - v) \cdot \left(\frac{\gamma}{h}\phi - \frac{\partial\phi}{\partial n}\right) - \int_{\Gamma_D} (v_n^D - v_n)\xi = 0$$

$$\int_{\Gamma_R} \mu(v_n^R - v_n)\left(\frac{\gamma}{h}\phi_n - \frac{\partial\phi}{\partial n}\cdot n\right) - \int_{\Gamma_R} (v_n^R - v_n)\xi = 0$$

$$\int_{\Gamma_P} \mu\pi_{n^\perp}(v^P - v)\cdot\left(\frac{\gamma}{h}\phi - \frac{\partial\phi}{\partial n}\right) = 0$$

Adding these terms we get

$$\begin{aligned}
a_\Omega(v,p;\phi\xi) - l_\Omega(\phi,\xi) =& -\int_{\Gamma_D}\frac{\gamma\mu}{h}v\cdot\phi + \int_{\Gamma_D}\mu\left(\frac{\partial v}{\partial n}\cdot\phi + v\cdot\frac{\partial\phi}{\partial n}\right) - \int_{\Gamma_D}(p\phi_n - v_n\xi)\\
&+\int_{\Gamma_D}\mu v^D\cdot\left(\frac{\gamma}{h}\phi - \frac{\partial\phi}{\partial n}\right) - \int_{\Gamma_D} v_n^D\xi\\
&-\int_{\Gamma_R}\frac{\gamma\mu}{h}v_n\phi_n + \int_{\Gamma_R}\mu\left(\frac{\partial v}{\partial n}\cdot n\phi_n + v_n\frac{\partial\phi}{\partial n}\cdot n\right) - \int_{\Gamma_R}(p\phi_n - v_n\xi)\\
&+\int_{\Gamma_R}\mu v_n^R\cdot\left(\frac{\gamma}{h}\phi_n - \frac{\partial\phi}{\partial n}\cdot n\right) - \int_{\Gamma_R} v_n^R\xi\\
&-\int_{\Gamma_P}\frac{\gamma\mu}{h}\pi_{n^\perp}v\cdot\pi_{n^\perp}\phi + \int_{\Gamma_P}\mu\left(\pi_{n^\perp}\frac{\partial v}{\partial n}\cdot\pi_{n^\perp}\phi + \pi_{n^\perp}v\cdot\pi_{n^\perp}\frac{\partial\phi}{\partial n}\right)\\
&+\int_{\Gamma_P}\mu\pi_{n^\perp}v^P\cdot\pi_{n^\perp}\left(\frac{\gamma}{h}\phi - \frac{\partial\phi}{\partial n}\right)\\
=& l_{\partial\Omega}(\phi,\xi) - a_{\partial\Omega}(v,p;\phi,\xi)
\end{aligned}$$

$\square$

Alternatively, we can write the system as

$$\begin{cases}
(v,p)\in V_h\times Q_h: \quad a(v,p;\phi\xi) + b(v,\xi) - b(\phi,p) = l_\Omega(\phi,\xi) + l_{\partial\Omega}(\phi,\xi) \quad \forall(\phi,\xi)\in V_h\times Q_h\\
a(v,p;\phi,\xi) := \int_\Omega \mu\nabla v:\nabla\phi + \lambda_R\int_{\Gamma_R}(v\cdot\phi - v_n\phi_n) + \int_{\Gamma_D\cup\Gamma_P}\frac{\gamma\mu}{h}v\cdot\phi - \int_{\Gamma_D\cup\Gamma_P}\mu\left(\frac{\partial v}{\partial n}\cdot n\phi + vn\cdot\frac{\partial\phi}{\partial n}\right)\\
\qquad +\int_{\Gamma_R}\frac{\gamma\mu}{h}v_n\phi_n - \int_{\Gamma_R}\mu\left(\frac{\partial v}{\partial n}\cdot n\phi_n + v_n\frac{\partial\phi}{\partial n}\cdot n\right) - \int_{\Gamma_P}\frac{\gamma\mu}{h}v_n\phi_n + \int_{\Gamma_P}\mu\left(\frac{\partial v}{\partial n}\cdot n\phi_n + v_n\frac{\partial\phi}{\partial n}\cdot n\right)\\
b(v,\xi) := \int_\Omega \mathrm{div}\, v\xi - \int_{\Gamma_D} v_n\xi
\end{cases}$$

$$(3.12)$$

### 3.4 Computation of boundary forces

Suppose $\psi\in\mathbb{R}^d$ is a vector field wich equals $\vec{d}\in\mathbb{R}^d$ on a given part $\Gamma\subset\partial\Omega$ of the boundary and vanishes on its complement. Then we get the sum of the forces on $\Gamma$ in direction $\vec{d}$ by means of the

integration by parts formula (3.7), supposed the weak solution is sufficiently smooth, as

$$\int_\Gamma n^\mathsf{T}\sigma\vec{d} = \int_{\partial\Omega} n^\mathsf{T}\sigma\psi = \int_\Omega \sigma : \nabla\psi + \int_\Omega \operatorname{div}\sigma\cdot\psi = \int_\Omega \sigma : \nabla\psi - \int_\Omega f\cdot\psi$$

$$= a_\Omega(v,p;\psi,0) - l_\Omega(\psi,0) - \lambda_R\int_{\Gamma_R}(v\cdot\psi - v_n\psi_n) + \int_{\Gamma_R}(g^R\cdot\psi - g_n^R\psi_n) - \int_{\Gamma_N}p^N\psi_n - \int_{\Gamma_P}p^P\psi_n$$

Supposing $\psi$ is a discrete vector field (in general we have to approximate it), for the strong implementation, we can retrieve the last expression by the parts of the matrix eliminated. Fo the weak implementation we have, since $\psi$ is now an admissible test function

$$\int_\Gamma n^\mathsf{T}\sigma\vec{d} = a_\Omega(v,p;\psi,0) - l_\Omega(\psi,0) - \lambda_R\int_{\Gamma_R}(v\cdot\psi - v_n\psi_n) + \int_{\Gamma_R}(g^R\cdot\psi - g_n^R\psi_n) - \int_{\Gamma_N}p^N\psi_n - \int_{\Gamma_P}p^P\psi_n$$

$$= l_{\partial\Omega}(\psi,0) - a_{\partial\Omega}(v,p;\psi,0) - \lambda_R\int_{\Gamma_R}(v\cdot\psi - v_n\psi_n) + \int_{\Gamma_R}(g^R\cdot\psi - g_n^R\psi_n) - \int_{\Gamma_N}p^N\psi_n - \int_{\Gamma_P}p^P\psi_n$$

$$= \int_{\Gamma_D}\mu v^D\cdot\frac{\gamma}{h}\psi + \int_{\Gamma_R}\mu v_n^R\cdot\frac{\gamma}{h}\psi_n + \int_{\Gamma_P}\mu\pi_{n\perp}v^P\cdot\pi_{n\perp}\frac{\gamma}{h}\psi$$

$$- \int_{\Gamma_D\cup\Gamma_P}\frac{\gamma\mu}{h}v\cdot\psi + \int_{\Gamma_D\cup\Gamma_P}\mu\frac{\partial v}{\partial n}\cdot\psi$$

$$- \int_{\Gamma_R}\frac{\gamma\mu}{h}v_n\psi_n - \int_{\Gamma_R}\mu\frac{\partial v}{\partial n}\cdot n\psi_n$$

$$+ \int_{\Gamma_P}\frac{\gamma\mu}{h}v_n\psi_n + \int_{\Gamma_P}\mu\frac{\partial v}{\partial n}\cdot n\psi_n$$

$$- \lambda_R\int_{\Gamma_R}(v\cdot\psi - v_n\psi_n) + \int_{\Gamma_R}(g^R\cdot\psi - g_n^R\psi_n) - \int_{\Gamma_N}p^N\psi_n - \int_{\Gamma_P}p^P\psi_n$$

$$\int_\Gamma n^\mathsf{T}\sigma\vec{d} = \begin{cases} \int_\Gamma\left(\mu\frac{\partial v}{\partial n} - pn\right)\cdot d + \int_\Gamma\frac{\gamma\mu}{h}(v^D - v)\cdot d & \Gamma\subset\Gamma_D \\ \int_\Gamma\mu\left(\frac{\partial v}{\partial n}\cdot n\phi_n\right) - \int_\Gamma p\phi_n + \int_\Gamma\mu(v_n^R - v_n)\cdot\left(\frac{\gamma}{h}\phi_n - \frac{\partial\phi}{\partial n}\cdot n\right) & \Gamma\subset\Gamma_R \end{cases} \tag{3.13}$$

### 3.5   Pressure mean

If no boundary conditions is of Neumann or Pressure type, the pressure is only determined up to a constant. In order to impose the zero mean on the pressure, let $C$ the matrix of size $(1,nc)$ with $C_{1j} = \int_\Omega \xi_j$.

### 3.6   Iterative solution

#### 3.6.1   ABCD

A general two-by tow system can be factorized, $A$ is inverible as well as $S = D - CA^{-1}B$ as

$$\mathcal{A} = \begin{bmatrix} A & B \\ C & D \end{bmatrix} = \begin{bmatrix} I & 0 \\ CA^{-1} & I \end{bmatrix}\begin{bmatrix} A & B \\ 0 & S \end{bmatrix}$$

so

$$\mathcal{A}\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} b \\ c \end{bmatrix} \quad \Leftrightarrow \quad \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} A^{-1}(b - By) \\ S^{-1}(c - CA^{-1}b) \end{bmatrix}$$

since

$$\begin{bmatrix} A & B \\ 0 & S \end{bmatrix}\begin{bmatrix} A^{-1}(b - By) \\ S^{-1}(c - CA^{-1}b) \end{bmatrix} = \begin{bmatrix} b \\ c - CA^{-1}b \end{bmatrix} \quad \Rightarrow \quad \begin{bmatrix} I & 0 \\ CA^{-1} & I \end{bmatrix}\begin{bmatrix} b \\ c - CA^{-1}b \end{bmatrix} = \begin{bmatrix} b \\ c \end{bmatrix}$$

### 3.6.2 Stokes

$$\mathcal{A} = \begin{bmatrix} A & -B^{\mathsf{T}} \\ B & 0 \end{bmatrix} = \begin{bmatrix} I & 0 \\ BA^{-1} & I \end{bmatrix}\begin{bmatrix} A & -B^{\mathsf{T}} \\ 0 & S \end{bmatrix} \qquad \mathcal{A}^{-1} = \begin{bmatrix} A^{-1} & A^{-1}B^{\mathsf{T}}S^{-1} \\ 0 & S^{-1} \end{bmatrix}\begin{bmatrix} I & 0 \\ -BA^{-1} & I \end{bmatrix} \qquad S = BA^{-1}B^{\mathsf{T}}$$

This means that the solution of $\mathcal{A}(x_v, x_p) = (y_v, y_p)$ is given by

$$\begin{cases} Ax_v' = y_v \\ Sx_p' = y_p - Bx_v' \\ Ax_v = y_v + B^{\mathsf{T}}x_p' \end{cases}$$

### 3.6.3 What is the effect of replacing $S$ by $K$?

Let

$$\mathcal{M}^{-1} := \begin{bmatrix} A^{-1} & A^{-1}B^{\mathsf{T}}K^{-1} \\ 0 & K^{-1} \end{bmatrix}\begin{bmatrix} I & 0 \\ -BA^{-1} & I \end{bmatrix} = \begin{bmatrix} A^{-1}(I - B^{\mathsf{T}}K^{-1}BA^{-1}) & A^{-1}B^{\mathsf{T}}K^{-1} \\ -K^{-1}BA^{-1} & K^{-1} \end{bmatrix}$$

Then we have

$$\mathcal{A}\mathcal{M}^{-1} = \begin{bmatrix} I & 0 \\ (I - SK^{-1})BA^{-1} & SK^{-1} \end{bmatrix}$$

### 3.6.4 What is the effect of replacing $A$ by $X$?

Let

$$\mathcal{M} := \begin{bmatrix} X & -B^{\mathsf{T}} \\ B & 0 \end{bmatrix}$$

Then with $T = BX^{-1}B^{\mathsf{T}}$

$$\begin{aligned}
\mathcal{M}^{-1}\mathcal{A} &= \begin{bmatrix} X^{-1} & X^{-1}B^{\mathsf{T}}T^{-1} \\ 0 & T^{-1} \end{bmatrix}\begin{bmatrix} I & 0 \\ -BX^{-1} & I \end{bmatrix}\begin{bmatrix} A & -B^{\mathsf{T}} \\ B & 0 \end{bmatrix} \\
&= \begin{bmatrix} X^{-1} & X^{-1}B^{\mathsf{T}}T^{-1} \\ 0 & T^{-1} \end{bmatrix}\begin{bmatrix} A & -B^{\mathsf{T}} \\ B(I - X^{-1}A) & T \end{bmatrix} \\
&= \begin{bmatrix} X^{-1}(A + B^{\mathsf{T}}T^{-1}B(I - X^{-1}A)) & 0 \\ T^{-1}B(I - X^{-1}A) & I \end{bmatrix}
\end{aligned}$$

We have

$$BX^{-1}(A + B^{\mathsf{T}}T^{-1}B(I - X^{-1}A)) = B$$

### 3.6.5 HSS preconditioning

$$\mathcal{A} = \mathcal{B} + \mathcal{C}, \quad \mathcal{B} = \begin{bmatrix} A & 0 \\ 0 & 0 \end{bmatrix}, \quad \mathcal{C} = \begin{bmatrix} 0 & -B^\mathsf{T} \\ B & 0 \end{bmatrix}, \quad \mathcal{M} = (\mathcal{B} + \alpha\mathfrak{I})(\mathcal{C} + \alpha\mathfrak{I})$$

Then

$$\begin{bmatrix} A + \alpha I & 0 \\ 0 & \alpha I \end{bmatrix} \begin{bmatrix} \alpha I & -B^\mathsf{T} \\ B & \alpha I \end{bmatrix} \begin{bmatrix} w \\ q \end{bmatrix} = \begin{bmatrix} v \\ p \end{bmatrix} \tag{3.14}$$

$$\begin{cases} (A + \alpha I)w' = v \\ q' = \dfrac{1}{\alpha}p \\ (\alpha I + \alpha^{-1}BB^\mathsf{T})p = q' - \alpha^{-1}Bw' \\ \alpha w = w' + B^\mathsf{T}p \end{cases}$$

### 3.6.6 System with constraint on pressure

We have to solve (**??**) with

$$\mathcal{A} = \begin{bmatrix} A & -B^\mathsf{T} & 0 \\ B & 0 & C^\mathsf{T} \\ 0 & C & 0 \end{bmatrix} = \begin{bmatrix} I & 0 & 0 \\ BA^{-1} & I & 0 \\ 0 & CS^{-1} & I \end{bmatrix} \begin{bmatrix} A & 0 & 0 \\ 0 & S & 0 \\ 0 & 0 & T \end{bmatrix} \begin{bmatrix} I & -A^{-1}B^\mathsf{T} & 0 \\ 0 & I & S^{-1}C^\mathsf{T} \\ 0 & 0 & I \end{bmatrix}$$

where $S = BA^{-1}B^\mathsf{T}$, $T = -CS^{-1}C^\mathsf{T}$. We have

$$\mathcal{A}^{-1} = \begin{bmatrix} I & A^{-1}B^\mathsf{T} & 0 \\ 0 & I & -S^{-1}C^\mathsf{T} \\ 0 & 0 & I \end{bmatrix} \begin{bmatrix} A^{-1} & 0 & 0 \\ 0 & S^{-1} & 0 \\ 0 & 0 & T^{-1} \end{bmatrix} \begin{bmatrix} I & 0 & 0 \\ -BA^{-1} & I & 0 \\ 0 & -CS^{-1} & I \end{bmatrix}$$

We construct our preconditioner by approximations of A, S, and T. The preconditioner $(y_{v,p}, y_\lambda) \rightarrow (x_v, x_p, x_\lambda)$ has the steps

$$\begin{cases} Ax'_v = y_v \\ Sx'_p = y_p - Bx'_v \\ Tx_\lambda = y_\lambda - Cx'_p \\ Sx''_p = C^\mathsf{T}x_\lambda \\ x_p = x'_p - x''_p \\ Ax''_v = B^\mathsf{T}x_p \\ x_v = x'_v + x''_v \end{cases}$$

# 4 Beam problem

We consider the beam equation with standard boundary conditions. Let $\Gamma_C \subset \partial\Omega$, $\Gamma_S \subset \partial\Omega$, and $\Gamma_F \subset \partial\Omega$ be the points where the clamped, simply supported and fixed boundary conditions hold.

$$\frac{d^2}{dx^2}(EI\frac{d^2w}{dx^2})(x) = q(x) \quad \Omega = ]0; L[$$

$$\begin{cases} w(x) = \frac{dw}{dx}(x) = 0 & \text{on } \Gamma_C \quad \text{(clamped end)} \\[2mm] w(x) = \frac{d^2w}{dx^2}(x) = 0 & \text{on } \Gamma_S \quad \text{(simply supported end)} \\[2mm] \frac{d^2w}{dx^2}(x) = \frac{\alpha}{EI}, \ \frac{d^3w}{dx^3}(x) = \frac{\beta}{EI} & \text{on } \Gamma_F \quad \text{(free end with forces)} \end{cases} \tag{4.1}$$

## 4.1 Weak formulation

Let

$$H_0^2(\Omega) := \left\{ v \in H^2(\Omega) \ \middle|\ v(x) = \frac{dv}{dx}(x) = 0 \,, x \in \partial\Omega \right\}. \tag{4.2}$$

be the standard Sobolov space. For $H_0^2(\Omega) \subset V \subset H^2(\Omega)$ we consider

$$w \in V : \quad \int_\Omega EI\frac{d^2w}{dx^2}\frac{d^2v}{dx^2} = \int_\Omega qv \quad \forall v \in V. \tag{4.3}$$

Integration by parts

$$\int_\Omega EI\frac{d^2w}{dx^2}\frac{d^2v}{dx^2} = -\int_\Omega \frac{d}{dx}(EI\frac{d^2w}{dx^2})\frac{dv}{dx} + \left[ EI\frac{d^2w}{dx^2}\frac{dv}{dx} \right]_0^L$$

$$= \int_\Omega \frac{d^2}{dx^2}(EI\frac{d^2w}{dx^2})v + \left[ EI\frac{d^2w}{dx^2}\frac{dv}{dx} \right]_0^L - \left[ EI\frac{d^3w}{dx^3}v \right]_0^L$$

shows that a regular solution of (4.3) satisfies

$$\left[ EI\frac{d^2w}{dx^2}\frac{dv}{dx} \right]_0^L - \left[ EI\frac{d^3w}{dx^3}v \right]_0^L = \int_\Omega \left( q - \frac{d^2}{dx^2}(EI\frac{d^2w}{dx^2}) \right) v.$$

Since $V$ is dense in $L^2(\Omega)$, choosing first $v \in H_0^2(\Omega)$ satisfies $\frac{d^2}{dx^2}(EI\frac{d^2w}{dx^2}) = q$ almost everywhere. Then we have

$$\left[ EI\frac{d^2w}{dx^2}\frac{dv}{dx} \right]_0^L - \left[ EI\frac{d^3w}{dx^3}v \right]_0^L = 0 \quad \forall v \in V.$$

So, depending on the choice of $V$ we have the following possible boundary conditions for $x \in \partial\Omega$:

| | constraint in V | strong b.c. | implied b.c. |
|---|---|---|---|
| 1 | $v(x) = \frac{dv}{dx}(x) = 0$ | $w(x) = \frac{dw}{dx}(x) = 0$ | – |
| 2 | $v(x) = 0$ | $w(x) = 0$ | $\frac{d^2w}{dx^2} = 0$ |
| 3 | $\frac{dv}{dx}(x) = 0$ | $\frac{dw}{dx}(x) = 0$ | $\frac{d^3w}{dx^3} = 0$ |
| 4 | – | – | $\frac{d^2w}{dx^2} = \frac{d^3w}{dx^3} = 0$ |

15

The third case does not seem to be of practical use. It follows that an appropriate choice for our problem is:

$$V := \left\{ v \in H^2(\Omega) \,\middle|\, v(x_c) = \frac{dv}{dx}(x_c) = 0, \quad v(x_s) = 0, \quad x_c \in \Gamma_C, x_s \in \Gamma_S \right\} \tag{4.4}$$

In order to take into account the non-homogenuous boundary condition we define

$$w \in V: \quad \int_\Omega EI \frac{d^2w}{dx^2} \frac{d^2v}{dx^2} = \int_\Omega qv + \delta_{x_f=L}(\alpha \frac{dv}{dx}(L) - \beta v(L)) - \delta_{x_f=0}(\alpha \frac{dv}{dx}(0) - \beta v(0)) \quad \forall v \in V. \tag{4.5}$$

For $a \in L^2(\Omega)$

**Proposition 4.1.** *(4.3) has a unique solution if $\Gamma_C \neq \emptyset$ and the solution satisfies a weak version of (4.1).*

*Proof.* Existence and uniqueness follow from the Lax-Milgram lemma and Poincaré's inequality, applied to the function and its derivative. □

## 4.2 Low order approximation

We use the spaces of quadratic B-splines on a one-dimensional mesh $h : 0 = x_0 < x_1 < \cdots < x_N = L$. The B-splines are the subspace of quadratic finite elements, which are of class $C^1$, allowing us to use a standard finite element basis.

Let $(\phi_i)_{0 \leqslant i \leqslant N}$ be the canonical bases $\mathcal{P}_h^1$ and $\psi_i(x) := \frac{(x-x_{i-1})(x_i-x)}{2h_i^2}$, $1 \leqslant i \leqslant N$. In addition let $h_i := x_i - x_{i-1}$ and $x_{i-\frac{1}{2}} := \frac{x_{i-1}+x_i}{2}$, $1 \leqslant i \leqslant N$.

Noticing that, with $u'$ the piecewise derivative of $u \in \mathcal{P}_h^2$, we have

$$u \in C^1(\Omega) \quad \Leftrightarrow \quad \sum_{j=1}^{N-1} [u'] \phi(x_j) = 0 \quad \forall 1 \leqslant i < N, \quad \Leftrightarrow \quad \int_\Omega \left( u'\phi_i' + u''\phi_i \right) = 0 \quad \forall 1 \leqslant i < N. \tag{4.6}$$

Taking into account the condition on the normal derivative in (4.4), we define

$$V_h := \left\{ v \in \mathcal{P}_h^2 \,\middle|\, \int_\Omega \left( v'\phi_i' + v''\phi_i \right) = 0 \quad \forall 0 \leqslant i \leqslant N \right\} \cap H_0^1(\Omega). \tag{4.7}$$

and the discrete problem is

$$w \in V_h: \quad \int_\Omega EI w''v'' = \int_\Omega fv \quad \forall v \in V_h. \tag{4.8}$$

Since (4.8) is the optimality condition of energy minimisation, we take into account the constraints defining $V_h$ by a Lagrange multiplier and use the splitting of the space $V_h = V_h^1 \times V_h^2$ and $\Lambda_1 = V_h^1$

$$w = \sum_{j=0}^{N} \vec{w}_{1j}\phi_j + \sum_{j=1}^{N} \vec{w}_{2j}\psi_j =: w_1 + w_2, \quad \lambda := \sum_{j=0}^{N} \vec{\lambda}_j \phi_j. \tag{4.9}$$

We notice that

$$v_1'' = 0 \quad \forall v_1 \in V_h^1, \quad \int_\Omega v_2'\mu' = 0 \quad \forall v_2 \in V_h^2, \mu \in \Lambda_h, \tag{4.10}$$

so we arrive at

$$(w_1, w_2, \lambda) \in V_h^1 \times V_h^2 \times \Lambda_h$$

$$
\begin{cases}
\displaystyle\int_\Omega v_1'\lambda' = \int_\Omega f v_1 & \forall v_1 \in V_h^1, \\[2mm]
\displaystyle\int_\Omega EI w_2'' v_2'' + \int_\Omega v_2''\lambda = \int_\Omega f v_2 & \forall v_2 \in V_h^2, \\[2mm]
\displaystyle\int_\Omega \left(w_1'\mu' + w_2''\mu\right) = 0 & \forall \mu \in \Lambda_h.
\end{cases}
$$

The boundary conditions correspond to

| | b.c. | constraint in $V_h^1$ | constraint in $\Lambda_h$ |
|---|---|---|---|
| 1 | clamped | $v(x) = 0$ | – |
| 2 | simply supp. | $v(x) = 0$ | $\mu(x) = 0$ |
| 3 | unknown | – | – |
| 4 | forces | – | $\mu(x) = 0$ |

We use again Lagrange multipliers for these constraints. Letting $\mathcal{C}_V, \mathcal{C}_\Lambda \subset \{0, L\}$

$$D_{ij} = \int_\Omega EI\psi_i''\psi_j'', \quad A_{ij} = \int_\Omega \mu_i'\phi_j', \quad B_{ij} = \int_\Omega \mu_i\psi_j'', \quad C_{1,ij} = \phi_j(x_i)\, x_i \in \mathcal{C}_V, \quad C_{2,ij} = \mu_j(x_i)\, x_i \in \mathcal{C}_\Lambda.$$

and

$$a_i := l(\phi_i), \quad b_i := l(\psi_i), \quad c_i = \mu(x_i)\frac{\partial u^D}{\partial n}(x_i), \quad d_i = \phi(x_i)u^D(x_i)$$

the discrete system reads

$$
\begin{bmatrix}
0 & 0 & A & C_1{}^{\mathsf{T}} & 0 \\
0 & D & B^{\mathsf{T}} & 0 & 0 \\
A & B & 0 & 0 & C_2{}^{\mathsf{T}} \\
C_1 & 0 & 0 & 0 & 0 \\
0 & 0 & C_2 & 0 & 0
\end{bmatrix}
\begin{bmatrix}
\vec{v_1} \\ \vec{v_2} \\ \vec{\lambda} \\ \alpha \\ \beta
\end{bmatrix}
=
\begin{bmatrix}
a \\ b \\ c \\ d \\ 0
\end{bmatrix},
\tag{4.11}
$$

Since D is a regular diagonal matrix we can eliminate $\vec{v_2}$:

$$\vec{v_2} = D^{-1}(b - B^{\mathsf{T}}\vec{\lambda}).$$

$$
\begin{bmatrix}
0 & A & C_1{}^{\mathsf{T}} & 0 \\
A & -X & 0 & C_2{}^{\mathsf{T}} \\
C_1 & 0 & 0 & 0 \\
0 & C_2 & 0 & 0
\end{bmatrix}
\begin{bmatrix}
\vec{v_1} \\ \vec{\lambda} \\ \alpha \\ \beta
\end{bmatrix}
=
\begin{bmatrix}
a \\ c - BD^{-1}b \\ d \\ 0
\end{bmatrix}, \quad X := BD^{-1}B^{\mathsf{T}}
$$

We have

$$\psi_i'(x) = \frac{(x_{i-\frac{1}{2}} - x)}{h_i^2}, \quad \psi_i''(x) = \frac{-1}{h_i^2},$$

$$B_{ii} = \int_{x_{i-1}}^{x_i} \phi_i\psi_i'' = \frac{-1}{2h_i}, \quad B_{i,i+1} = \frac{-1}{2h_{i+1}}, \quad D_{ii} = \frac{EI_i}{h_i^3}$$

$$
\begin{cases}
X_{i,i-1} = \dfrac{h_i}{4EI_i} \\[2mm]
X_{i,i} = \dfrac{h_i}{4EI_i} + \dfrac{h_{i+1}}{4EI_{i+1}} \\[2mm]
X_{i,i+1} = \dfrac{h_{i+1}}{4EI_{i+1}}
\end{cases}
$$

17

$$\begin{bmatrix} A & B^\mathsf{T} \\ B & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix} \quad \Rightarrow \quad BB^\mathsf{T}y = Ba - BAx \quad \Rightarrow \quad (I-C)Ax = (I-C)a, \quad C := B^\mathsf{T}\left(BB^\mathsf{T}\right)^{-1}B$$

$$\Rightarrow \quad ((I-C)A + C)\,x = (I-C)a + B^\mathsf{T}\left(BB^\mathsf{T}\right)^{-1}b$$

# A Linear Algebra

## A.1 Saddle-point systems

Let $(X, \langle \cdot, \cdot \rangle)$ be a Hilbert space and

$$X_c := \{x \in X \mid Bx = c\}.$$

We consider the minimization problem with symmetric $A \in \mathcal{L}(X, X)$ and $B \in \mathcal{L}(X, Y)$

$$\inf \left\{ \frac{1}{2} \langle Ax, x \rangle - \langle b, x \rangle \;\middle|\; x \in X_c \right\} \tag{A.1}$$

With the Lagrangian $L(x, y) := J(x) + \langle Bx - c, y \rangle$, $J(x) := \frac{1}{2} \langle Ax, x \rangle - \langle b, x \rangle$ we get the saddle-point system

$$\begin{bmatrix} A & B^\mathsf{T} \\ B & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} b \\ c \end{bmatrix} \tag{A.2}$$

### A.1.1 Solution by the kernel method

Since $X_0 = \ker B$, we may write $X = X_0 \oplus X_1$ and split the vectors and matrices accordingly.

$$\begin{bmatrix} A_{00} & A_{01} & B_0^\mathsf{T} \\ A_{10} & A_{11} & B_1^\mathsf{T} \\ B_0 & B_1 & 0 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ y \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \\ c \end{bmatrix} \tag{A.3}$$

But by definition we have $B_0 = 0$ and we can rewrite (A.4) as

$$\begin{bmatrix} A_{00} & A_{01} & 0 \\ 0 & B_1 & 0 \\ A_{10} & A_{11} & B_1^\mathsf{T} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ y \end{bmatrix} = \begin{bmatrix} b_0 \\ c \\ b_1 \end{bmatrix} \tag{A.4}$$

The system is well-posed, since $B_1$ is invertible.

### A.1.2 Penalty and augmented Lagrangian method

The penalty method is the unconstrained minimization of $J_\varepsilon(x) := J(x) + \frac{1}{2\varepsilon} \|Bx - c\|_{M^{-1}}^2$ and leads to the penalty system

$$(A + \frac{1}{\varepsilon} B^\mathsf{T} M^{-1} B)x = b + \frac{1}{\varepsilon} B^\mathsf{T} M^{-1} c. \tag{A.5}$$

It is equivalent to the modified system

$$\begin{bmatrix} A & B^\mathsf{T} \\ B & -\varepsilon M \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} b \\ c \end{bmatrix}$$

The augmented Lagrangian method uses the Lagrangian $L_r(x, y) := J(x) + \frac{r}{2} \|Bx - c\|_{M^{-1}}^2 + \langle Bx - c, y \rangle$ and leads to the saddle-point system

$$\begin{bmatrix} (A + rB^\mathsf{T} M^{-1} B) & B^\mathsf{T} \\ B & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} b + rB^\mathsf{T} M^{-1} c \\ c \end{bmatrix}$$

The first diagonal equals the one of the penalty method for $r = 1/\varepsilon$, but the method is always conforming in the sense that the constraints are satisfied exactly, which allows for arbitrary choice of $r$. If we combine with withe kernel method, we get

$$
\begin{bmatrix} A_{00} & A_{01} & 0 \\ 0 & B_1 & 0 \\ A_{10} & (A_{11} + rB_1^T M^{-1} B_1) & B_1^T \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ y \end{bmatrix} = \begin{bmatrix} b_0 \\ c \\ b_1 + rB_1^T M^{-1} c \end{bmatrix},
$$

which shows that only the computation of the multiplier $y$ is changed, without changing it.

A different situation arrives in the special case $B = \begin{bmatrix} 0 & A_{11} \end{bmatrix}$ and $M = A_{11}$. Then we have

$$
\begin{bmatrix} A_{00} & A_{01} & 0 \\ A_{10} & (1+r)A_{11} & A_{11} \\ 0 & A_{11} & 0 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ y \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 + rA_{11}c \\ c \end{bmatrix}
$$

# B  Finite elements on simplices

## B.1  Simplices

We consider an arbitrary non-degenerate simplex $K = (x_0, x_1, \ldots, x_d)$. The volume of K is given by

$$|K| = \frac{1}{d!} \det(x_1 - x_0, \ldots, x_d - x_0) = \frac{1}{d!} \det(1, x_0, x_1 \ldots, x_d) \quad 1 = (1, \ldots, 1)^{\mathsf{T}1}. \tag{B.1}$$

The d+1 sides $S_k$ (co-dimension one, $d-1$-simplices or facets) are defined by $S_k = (x_0, \ldots, \cancel{x_k}, \ldots, x_d)$. The height is $h_k = |P_{S_k} x_k - x_k|$, where $P_S$ is the orthogonal projection on the hyperplane associated to $S_k$. We have $P_{S_k} x_k = x_k + h_k \vec{n}[k]$ and $S_k = \left\{ x \in \mathbb{R}^d \mid \vec{n}[k]^{\mathsf{T}} x = h_k \right\}$ and

$$0 = \int_K \operatorname{div}(\vec{c}) = \sum_{i=0}^{d} \int_{S_i} \vec{c} \cdot \vec{n}[i] = \vec{c} \cdot \sum_{i=0}^{d} |S_i| \vec{n}[i] \quad \Rightarrow \quad \sum_{i=0}^{d} |S_i| \vec{n}[i] = 0$$

$$d |K| = \int_K \operatorname{div}(x) = \sum_{i=0}^{d} \int_{S_i} x \cdot \vec{n}[i] = \sum_{i=0}^{d} |S_i| h_i$$

> **Height formula**
>
> $$h_k = d \frac{|K|}{|S_k|}$$

## B.2  Barycentric coordinates

The barycentric coordinate of a point $x \in \mathbb{R}^d$ give the coefficients in the affine combination of $x = \sum_{i=0}^{d} \lambda_i x_i$ ($\sum_{i=0}^{d} \lambda_i = 1$) and can be expressed by means of the outer unit normal $\vec{n}[i]$ of $S_i$ or the signed distance $d^s$ as

$$\lambda_i(x) = \frac{\vec{n}[i]^{\mathsf{T}}(x_j - x)}{\vec{n}[i]^{\mathsf{T}}(x_j - x_i)} \quad (j \neq i), \qquad \lambda_i(x) = \frac{d^s(x, H)}{h_i}. \tag{B.2}$$

Any polynomial in the barycentric coordinates can be integrated exactly. For $\alpha \in \mathbb{N}_0^{d+1}$ we let $\alpha! = \prod_{i=0}^{d} \alpha_i!$, $|\alpha| = \sum_{i=0}^{d} \alpha_i$, and $\lambda^\alpha = \prod_{i=0}^{d} \lambda_i^{\alpha_i}$

> **Integration on K**
>
> $$\int_K \lambda^\alpha = |K| \frac{d! \alpha!}{(|\alpha| + d)!} \tag{B.3}$$

see [**EisenbergMalvern73**], [**VermolenSegal18**].

---

[1] https://en.wikipedia.org/wiki/Simplex#Volume

## B.3 Finite elements

We consider a family $\mathcal{H}$ of regular simplicial meshes $h$ on a polyhedral domain $\Omega \subset \mathbb{R}^d$. The set of simplices of $h \in \mathcal{H}$ is denoted by $\mathcal{K}_h$, and its $d-1$-dimensional sides by $\mathcal{S}_h$, divided into interior and boundary sides $\mathcal{S}_h^{int}$ and $\mathcal{S}_h^\partial$, respectively. The set of $d+1$ sides of $K \in \mathcal{K}_h$ is $\mathcal{S}_h(K)$. To any side $S \in \mathcal{S}_h$ we associate a unit normal vector $n_S$, which coincides with the unit outward normal vector $n_{\partial\Omega}$ if $S \in \mathcal{S}_h^\partial$.

For $K \in \mathcal{K}_h$ and $S \in \mathcal{S}_h$, or $S \in \mathcal{S}_h(K)$ we denote

$$
\begin{array}{lll}
x_K : & \text{barycenter of } K & \qquad x_S : \quad \text{barycenter of } S \\
x_S^K : & \text{vertex opposite to } S \text{ in } K & \qquad h_S^K : \quad \text{distance of } x_S^K \text{ to } S \\
\sigma_S^K := \begin{cases} +1 & \text{if } n_S = n_K, \\ -1 & \text{if } n_S = -n_K. \end{cases} & & \qquad \lambda_S^K : \quad \text{barycentric coordinates of } K
\end{array}
$$

For $k \in \mathbb{N}_0$ we denote by $\mathcal{C}_h^k(\Omega)$ the space of piecewise $k$-times differential functions with respect to $\mathcal{K}_h$. The subspace of piecewise polynomial functions of order $k \in \mathbb{N}_0$ in $C_h^k(\Omega)$ is denoted by $\mathcal{D}_h^k(\Omega)$ and the $L^2(\Omega)$-projection by $\pi_h^k : L^2(\Omega) \to \mathcal{D}_h^k(\Omega)$.

### B.3.1 $\mathcal{P}_h^1(\Omega)$

We have $\mathcal{P}_h^1(\Omega) = \mathcal{D}_h^1(\Omega) \cap C(\overline{\Omega})$, but the FEM definition also provides a basis. The restrictions of the basis functions of $\mathcal{P}_h^1(\Omega)$ to the simplex $K$ are the barycentric coordinates $\lambda_S^K$ associated to the node opposite to $S$ in $K$.

For the computation of matrices we use (B.3), for example for $i, j \in [\![0, d]\!]$

$$
\int_K \lambda_i \lambda_j = |K| \frac{d! \alpha!}{(|\alpha| + d)!} \quad \text{with} \quad \begin{cases} \alpha = (1, 1, 0, \cdots, 0) & (i \neq j) \\ \alpha = (2, 0, \cdots, 0) & (i = j) \end{cases}
$$

so

$$
\int_K \lambda_i \lambda_j = \frac{|K|}{(d+2)(d+1)} (1 + \delta_{ij}) \tag{B.5}
$$

More generally, we have for $i_l \in [\![0, d]\!]$ with $1 \leqslant l \leqslant k$

$$
\int_K \lambda_{i_1} \cdots \lambda_{i_k} = \frac{|K| \alpha!}{(d+k) \cdots (d+1)}, \qquad \alpha_l = \# \{ j \in [\![0, d]\!] \mid i_j = l \}, \quad 1 \leqslant l \leqslant k. \tag{B.6}
$$

### B.3.2 $\mathcal{CR}_h^1(\Omega)$

$$\mathcal{CR}_h^k(\Omega) := \left\{ q \in \mathcal{D}_h^k(\Omega) \, \middle| \, \int_S [q]\, p = 0 \; \forall S \in \mathcal{S}_h^{int}, \forall p \in P^{k-1}(S) \right\}. \tag{B.7}$$

Denote in addition the basis of $\mathcal{CR}_h^1(\Omega)$ by $\psi_S$, we have

> **Formulae for $\mathcal{CR}_h^1$**
>
> $$\psi_S\big|_K = 1 - d\lambda_S^K, \quad \nabla \psi_S\big|_K = \frac{|S|\sigma_S^K}{|K|} n_S, \quad \frac{1}{|K|} \int_K \psi_S = \frac{1}{d+1}. \tag{B.8}$$

### B.3.3 $\mathcal{RT}_h^0(\Omega)$

The Raviart-Thomas space for $k \geqslant 0$ is given by

$$\mathcal{RT}_h^k(\Omega) := \left\{ v \in D_h^k(\Omega, \mathbb{R}^d) \oplus X_h^k \, \middle| \, \int_S [v_n]\, p = 0 \; \forall S \in \mathcal{S}_h^{int}, \forall p \in P^k(S) \right\} \tag{B.9}$$

where $X_h^k := \left\{ xp \mid p\big|_K \in P_{hom}^k(K) \; \forall K \in \mathcal{K}_h \right\}$ with $P_{hom}^k(K)$ the space of $k$-th order homogenous polynomials.

Then the Raviart-Thomas basis function of lowest order is given by

> **Formulae for $\mathcal{RT}^0$**
>
> $$\Phi_S\big|_K := \sigma_S^K \frac{x - x_S^K}{h_S^K}, \quad \int_K \operatorname{div} \Phi_S\big|_K = \sigma_S^K \frac{d|K|}{h_S^K} = \sigma_S^K |S|, \quad \frac{1}{|K|} \int_K \Phi_S = \sigma_S^K \frac{x_K - x_S^K}{h_S^K}. \tag{B.10}$$

For the `pyhon` implementation of the projection on $\mathcal{D}_h^0(\Omega, \mathbb{R}^d)$ we have with the height formula

$$\pi_h(\vec{v})\big|_K = \sum_{i=1}^d v_i \frac{1}{|K|} \int_K \Phi_i(x) = \sum_{i=1}^d v_i \sigma_i^K (x_K - x_{S_i}) \frac{|S_i|}{d\,|K|}$$

The `pyhon` implementation reads

### B.3.4 Moving a point to the boundary

Let $K$ be a simplex and $x \in K = \operatorname{conv}\{a_i \mid 0 \leqslant i \leqslant d\}$ given, i.e.

$$x = \sum_{i=0}^d \lambda_i a_i = a_0 + \sum_{i=1}^d \lambda_i (a_i - a_0)$$

Given $\beta \in \mathbb{R}^d$ we wish to find $x_\beta \in \partial K$ such that

$$x_\beta = \sum_{i=0}^d \mu_i a_i, \quad x_\beta = x + \delta\beta, \quad \delta > 0. \tag{B.11}$$

The condition $x_\beta \in \partial K$ amounts to $0 \leqslant \mu_i \leqslant 1$, $\sum_{i=0}^d \mu_i = 1$, and $\delta$ to be maximal. We get the solution in two steps. First we find $b_i$ such that

$$\beta = \sum_{i=1}^d b_i(a_i - a_0),$$

which gives

$$\sum_{i=1}^d (\mu_i - \lambda_i - \delta b_i)(a_i - a_0) = 0 \quad \Rightarrow \quad \mu_i = \lambda_i + \delta b_i \quad \forall 1 \leqslant i \leqslant d.$$

Now $\delta$ has to be chosen, such that the point $x_\beta$ lies inside $K$, i.e.

$$\begin{cases} 0 \leqslant \lambda_i + \delta b_i \leqslant 1 \\ 0 \leqslant \sum_{i=1}^d (\lambda_i + \delta b_i) \leqslant 1 \end{cases} \Leftrightarrow \begin{cases} -\lambda_i \leqslant \delta b_i \leqslant 1 - \lambda_i \quad \forall 1 \leqslant i \leqslant d, \\ \delta \sum_{i=1}^d b_i \leqslant \lambda_0 \end{cases}$$

**Lemma B.1.** *Let $0 \leqslant \lambda_i \leqslant 1$. Then the solution of*

$$\max \left\{ \delta \, \middle| \, -\lambda_i \leqslant \delta b_i \leqslant 1 - \lambda_i \quad \forall 1 \leqslant i \leqslant d, \quad \delta \sum_{i=1}^d b_i \leqslant \lambda_0 \right\} \tag{B.12}$$

*is*

$$\delta = \min \left\{ \min \left\{ \frac{1 - \lambda_i}{b_i} \, \middle| \, b_i > 0 \right\}, \min \left\{ \frac{-\lambda_i}{b_i} \, \middle| \, b_i < 0 \right\}, \frac{\lambda_0}{\sum_{i=1}^d b_i} \right\} \quad \textit{if} \quad \sum_{i=1}^d b_i > 0 \tag{B.13}$$

*Proof.* For $b_i > 0$ we have $\delta \leqslant \frac{1 - \lambda_i}{b_i}$, so $0 \leqslant \delta b_i + \lambda_i \leqslant 1$.
For $b_i < 0$ we have $\delta \leqslant \frac{-\lambda_i}{b_i}$, so $0 \leqslant \lambda_i + \delta b_i \leqslant \lambda_i \leqslant 1$. $\qquad \square$

# C  Discreization of the transport equation

For $k \in \mathbb{N}_0$ we denote by $\mathcal{C}_h^k(\Omega)$ the space of piecewise $k$-times differential functions with respect to $\mathcal{K}_h$, and piecewise differential operators $\nabla_h : \mathcal{C}_h^l(\Omega) \to \mathcal{C}_h^{l-1}(\Omega, \mathbb{R}^d)$ ($l \in \mathbb{N}$) by $\nabla_h q|_K := \nabla\left(q|_K\right)$ for $q \in \mathcal{C}_h^l(\Omega)$ and similarly for $\mathrm{div}_h : \mathcal{C}_h^l(\Omega, \mathbb{R}^d) \to \mathcal{C}_h^{l-1}(\Omega)$. We frequently use the piecewise Stokes formula

$$\int_\Omega (\nabla_h q)v + \int_\Omega q(\mathrm{div}_h v) = \int_{\mathcal{S}_h^{int}} [qv_n] + \int_{\mathcal{S}_h^\partial} qv_n, \tag{C.1}$$

where $\int_{\mathcal{S}_h} = \sum_{S \in \mathcal{S}_h} \int_S$ and $n$ in the sum stands for $n_S$.

The subspace of piecewise polynomial functions of order $k \in \mathbb{N}_0$ in $C_h^k(\Omega)$ is denoted by $\mathcal{D}_h^k(\Omega)$ and the $L^2(\Omega)$-projection by $\pi_h^k : L^2(\Omega) \to \mathcal{D}_h^k(\Omega)$.

Suppose $u$ satisfies

$$\mathrm{div}(\beta u) = f \quad \text{in } \Omega, \qquad \beta_n^-(u - u^D) = 0 \quad \text{on } \partial\Omega. \tag{C.2}$$

From the integration by parts formula

$$\int_\Omega \mathrm{div}(\beta u)v = -\int_\Omega \beta u \cdot \nabla v + \int_{\partial\Omega} \beta_n uv \tag{C.3}$$

it then follows that $u$ satisfies

$$a(u, v) = l(v) \quad \forall v \in V$$

with

$$a(u, v) := \int_\Omega \mathrm{div}(\beta u)v - \int_{\partial\Omega} \beta_n^- uv, \quad l(v) := \int_\Omega fv - \int_{\partial\Omega} \beta_n^- u^D v. \tag{C.4}$$

**Lemma C.1.**

$$a(u, u) = \int_\Omega \frac{\mathrm{div}(\beta)}{2} u^2 + \int_{\partial\Omega} \frac{|\beta_n|}{2} u^2. \tag{C.5}$$

*Proof.* We also have

$$a(u, v) = \frac{1}{2}\int_\Omega \mathrm{div}(\beta)uv + \frac{1}{2}\int_\Omega \mathrm{div}(\beta u)v + \frac{1}{2}\int_\Omega (\beta \cdot \nabla u)v - \int_{\partial\Omega} \beta_n^- uv$$

$$= \frac{1}{2}\int_\Omega \mathrm{div}(\beta)uv + \frac{1}{2}\int_\Omega ((\beta \cdot \nabla u)v - u(\beta \cdot \nabla v)) + \int_{\partial\Omega} (\frac{1}{2}\beta_n - \beta_n^-)uv$$

$$= \frac{1}{2}\int_\Omega \mathrm{div}(\beta)uv + \frac{1}{2}\int_\Omega ((\beta \cdot \nabla u)v - u(\beta \cdot \nabla v)) + \int_{\partial\Omega} \frac{|\beta_n|}{2}uv$$

such that the result follows with $v = u$. $\qquad\square$

## C.1  $\mathcal{D}_h^k(\Omega)$

Let

$$\begin{cases} a_h(u, v) := \int_\Omega \mathrm{div}_h(\beta u)v - \int_{\partial\Omega} \beta_n^- uv - \int_{\mathcal{S}_h} [u]\,\beta_S^\sharp(v) \\[2mm] \beta_S^\sharp(v) := \beta_{n_S}^- v^{in} + \beta_{n_S}^+ v^{ex} = \beta_{n_S}\{v\} - \frac{|\beta_{n_S}|}{2}[v] \end{cases} \tag{C.6}$$

**Lemma C.2.** *We have*

$$
\begin{cases}
a_h(u,v) = -\int_\Omega u(\beta \cdot \nabla_h v) + \int_{\partial\Omega} \beta_n^+ uv + \int_{\mathcal{S}_h} \beta_S^\flat(u)\,[v]\,, \\[2mm]
\beta_S^\flat(u) := \beta_{n_S}^+ u^{in} + \beta_{n_S}^- u^{ex} = -(-\beta_S)^\sharp(u) = \beta_{n_S}\{v\} + \dfrac{|\beta_{n_S}|}{2}[v]
\end{cases}
\tag{C.7}
$$

*and*

$$
a_h(u,v) = \frac{1}{2}\int_\Omega \left(\mathrm{div}_h(\beta u)v - u(\beta \cdot \nabla_h v)\right) + \int_{\partial\Omega} \frac{|\beta_n|}{2}uv + \int_{\mathcal{S}_h} \frac{|\beta_n|}{2}[u]\,[v] + \int_{\mathcal{S}_h} \frac{\beta_n}{2}\left(u^{ex}v^{in} - u^{in}v^{ex}\right)
\tag{C.8}
$$

*Proof.*

$$
\int_\Omega \mathrm{div}_h(\beta u)v = -\int_\Omega u(\beta \cdot \nabla_h v) + \int_{\partial\Omega} \beta_n uv + \int_{\mathcal{S}_h} \beta_{n_S}\,[uv]
$$

We get (C.7) with

$$
\beta_{n_S}\,[uv] - [u]\,\beta_S^\sharp(v) = \beta_{n_S}\left([u]\{v\} + \{u\}[v]\right) - [u]\,\beta_{n_S}\{v\} + \frac{|\beta_{n_S}|}{2}[u]\,[v]
$$

$$
= \beta_{n_S}\{u\}[v] + \frac{|\beta_{n_S}|}{2}[u]\,[v] = \beta_S^\flat(u)\,[v]\,.
$$

Finally for (C.8)

$$
\beta_S^\flat(u)\,[v] - [u]\,\beta_S^\sharp(v) = |\beta_n|\,[u]\,[v] + \beta_{n_S}\{u\}[v] - [u]\,\beta_{n_S}\{v\}
$$

$$
\beta_{n_S}\{u\}[v] - [u]\,\beta_{n_S}\{v\} = \frac{\beta_n}{2}\left(u^{ex}v^{in} - u^{in}v^{ex}\right)
$$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

**Corollary C.3.**

$$
a_h(u,u) = \int_\Omega \frac{\mathrm{div}_h(\beta)}{2}u^2 + \int_{\partial\Omega} \frac{|\beta_n|}{2}u^2 + \int_{\mathcal{S}_h} \frac{|\beta_{n_S}|}{2}[u]^2
\tag{C.9}
$$

*Proof.*

$$
2a_h(u,u) = \int_\Omega \mathrm{div}_h(\beta u)u - \int_{\partial\Omega} \beta_n^- uu - \int_{\mathcal{S}_h} \beta_S^\sharp(u)\,[u] - \int_\Omega u(\beta \cdot \nabla_h u) + \int_{\partial\Omega} \beta_n^+ uu + \int_{\mathcal{S}_h} [u]\,\beta_S^\flat(u)
$$

$$
= \int_\Omega \mathrm{div}_h(\beta)u^2 + \int_{\partial\Omega} |\beta_n|\,u^2 + \int_{\mathcal{S}_h} [u]\left(\beta_S^\flat(u) - \beta_S^\sharp(u)\right)
$$

$$
\beta_S^\flat(u) - \beta_S^\sharp(u) = \beta_{n_S}^+ u^{in} + \beta_{n_S}^- u^{ex} - \beta_{n_S}^- u^{in} - \beta_{n_S}^+ u^{ex} = |\beta_{n_S}|\,u^{in} - |\beta_{n_S}|\,u^{ex}
$$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

We suppose $\beta \in \mathcal{RT}_h^1$ with $\mathrm{div}\,\beta = 0$. Then $\beta \in D_h^0$ and we have

$$
\int_\Omega u(\beta \cdot \nabla_h v) = \int_\Omega \pi_h u(\beta \cdot \nabla_h v) = \int_{\partial\Omega} \beta_n(\pi_h u)v + \int_{\mathcal{S}_h} \beta_n\,[\pi_h u]\,v
$$

**Corollary C.4.** *For* $k = 0$ *the solution to*

$$u \in \mathcal{D}_h^0: \quad a_h(u, v) = l(v) \quad \forall v \in \mathcal{D}_h^0 \tag{C.10}$$

*satisfies monotonicity:* $l \geq 0$ *implies* $u \geq 0$

*Proof.* We write $u = u^+ + u^-$ and use $v = u^-$ in (C.10) such that

$$a(u^-, u^-) = a(u, u^-) - a(u^+, u^-) = l(u^-) - a(u^+, u^-) \leq -a(u^+, u^-).$$

and since with $x - |x| = 2x^-$ and $-x - |x| = -2x^+$

$$\int_{\mathcal{S}_h} \frac{|\beta_n|}{2} [u] [v] + \int_{\mathcal{S}_h} \frac{\beta_n}{2} \left( u^{ex} v^{in} - u^{in} v^{ex} \right) = \int_{\mathcal{S}_h} \frac{|\beta_n|}{2} \left( u^{in} v^{in} + u^{ex} v^{ex} \right) + \int_{\mathcal{S}_h} \left( \beta_n^- u^{ex} v^{in} - \beta^+ u^{in} v^{ex} \right) \tag{C.11}$$

$$a(u^+, u^-) = \int_{\partial\Omega} \frac{|\beta_n|}{2} u^+ u^- + \int_{\mathcal{S}_h} \frac{|\beta_n|}{2} [u^+] [u^-] + \int_{\mathcal{S}_h} \frac{\beta_n}{2} \left( u^{+ex} u^{-in} - u^{+in} u^{-ex} \right)$$

$$= \underbrace{\int_{\partial\Omega} \frac{|\beta_n|}{2} u^+ u^-}_{=0} + \underbrace{\int_{\mathcal{S}_h} \frac{|\beta_n|}{2} \left( u^{+in} u^{-in} + u^{+ex} u^{-ex} \right)}_{=0} + \underbrace{\int_{\mathcal{S}_h} \left( \beta_n^- u^{+ex} u^{-in} - \beta^+ u^{+in} u^{-ex} \right)}_{\geq 0}$$

Since $a(u, u)$ is norm on $\mathcal{D}_h^0$, we have $u^- = 0$, i.e. $u \geq 0$. □

## C.2 $\mathcal{D}_h^1(\Omega)$

We have for $\beta \in \mathcal{RT}_h^0$ with $\mathrm{div}\, \beta = 0$

$$\int_\Omega (\beta \cdot \nabla_h u) v = \int_\Omega (\beta \cdot \nabla_h u) \pi_h^0 v = \int_{\mathcal{S}_h^{int}} \beta_n [u \pi_h^0 v] + \int_{\partial\Omega} u \beta_n \pi_h^0 v$$

## C.3 $\mathcal{P}_h^1(\Omega)$

Let $K \in \mathcal{K}_h$, $\beta_K = \pi_K \beta$, $x_K$ be the barycenter of $K$ and $x_K^\sharp \in \partial K$ such that with $\delta_K \geq 0$

$$x_K^\sharp = x_K + \delta_K \beta_K \tag{C.12}$$

If we know $\vec{n}[i]^\mathsf{T} \beta_K$, we can compute $x_K^\sharp$ as follows.

$$\lambda_i(x_K^\sharp) = \lambda_i(x_K) + \delta_K \nabla \lambda_i^\mathsf{T} \beta_K = \frac{1}{d+1} - \delta_K \frac{\vec{n}[i]^\mathsf{T} \beta_K}{h_i} = \frac{1}{d+1} - \delta_K \frac{\vec{n}[i]^\mathsf{T} \beta_K |S_i|}{d |K|}$$

It follows that

$$\delta_K = \max \left\{ \frac{d |K|}{(d+1) |S_i| \left( \vec{n}[i]^\mathsf{T} \beta_K \right)^+} \,\middle|\, 0 \leq i \leq d \right\}. \tag{C.13}$$

The stabilized bilinear form is

$$a^{supg}(u, v) := \int_\Omega (\beta \cdot \nabla u) v - \int_{\partial\Omega} \beta_n^- u v + \int_\Omega \delta(\beta \cdot \nabla u)(\beta \cdot \nabla v) \tag{C.14}$$

Then we have

$$a^{supg}(u, v) =$$

# D Python implementation

We suppose to have a `class SimplexMesh` contaning the following elements

```
class SimplexMesh ():
  dimension , nnodes , ncells , nfaces
  simplices # np.array (( ncells , dimension +1))
  faces      # np.array (( nfaces , dimension ))
  points , pointsc , pointsf # np.array (( nnodes ,3)) , np.array (( ncells ,3)) , np.array ((
  normals , sigma   # np.array (( nfaces ,dimension )) , np.array (( ncells , dimension +1))
  dV              # np.array (( ncells ))
  bdrylabels     # dictionary (keys : colors , values : id 's of boundary faces )
```

The norm of the 'normals' $\widetilde{\vec{n}}$ is the measure of of the face

$$\widetilde{\vec{n}}[i] = |S_i|\,\vec{n}[i]$$

28