# Numerical Analysis of Partial Differential Equations

Alfio Quarteroni

MOX, Dipartimento di Matematica
Politecnico di Milano

Lecture Notes
A.Y. 2022-2023

# Advection–diffusion–reaction (ADR) equations
## Cfr [Q], Chap. 13

We consider the problem $\mathcal{L}u = f$ in $\Omega$, $u = 0$ on $\partial\Omega$, where:

1. $\mathcal{L}u = -\operatorname{div}(\mu\nabla u + \mathbf{b}u) + \sigma u$ \qquad (conservative form)
2. $\mathcal{L}u = -\operatorname{div}(\mu\nabla u) + \mathbf{b}\cdot\nabla u + \sigma u$ \qquad (non-conservative form)

Assumptions on coefficients as in Lecture 1 (see slide 8).

Weak formulation:

$$
\begin{cases}
\text{Find } u \in V = H_0^1(\Omega) \\
a(u, v) = F(v) \qquad \forall\, v \in V
\end{cases}
\tag{1}
$$

$$
F(v) = \int_\Omega f\, v
$$

$$
a(u, v) = \begin{cases}
\int_\Omega (\mu\nabla u + \mathbf{b}u)\cdot\nabla v + \int_\Omega \sigma\, u\, v & \text{conservative} \\
\int_\Omega \mu\nabla u \cdot \nabla v + \int_\Omega \mathbf{b}\cdot\nabla u\, v + \int_\Omega \sigma\, u\, v & \text{non-conservative}
\end{cases}
\tag{2}
$$

[Q] A. Quarteroni, *Numerical Models for Differential Problems*, 3rd Ed., Springer, 2018

# Lax-Milgram Lemma hypotheses

**Coercivity**

Sufficient conditions for coercivity:

1. Non-conservative case:  $\sigma - \frac{1}{2} \operatorname{div} \mathbf{b} \geq 0$ in $\Omega$
   (see Lecture 1, slides 22–23)

2. Conservative case:  $\sigma + \frac{1}{2} \operatorname{div} \mathbf{b} \geq 0$ in $\Omega$
   (prove it – similar proof)

In both cases: $a(u, u) \geq \mu_0 \|\nabla u\|^2 \quad \rightarrow \quad$ coercivity constant $\alpha \simeq \mu_0$

**Continuity**

In both cases, continuity constant: $M \simeq \|\mu\|_{L^\infty} + \|\mathbf{b}\|_{L^\infty} + \|\sigma\|_{L^2}$

(see Lecture 1, slide 21, for the non-conservative case. Similar proof for the conservative case – see book [Q])

---

[Q] A. Quarteroni, *Numerical Models for Differential Problems*, 3rd Ed., Springer, 2018

# FE error estimate

$$\|u - u_h\| \overset{\text{(Céa)}}{\leq} \frac{M}{\alpha} \inf_{v_h \in V_h} \|u - v_h\| \overset{\substack{\text{(interpolation} \\ \text{error estimate)}}}{\leq} C\frac{M}{\alpha} h^r |u|_{H^{r+1}(\Omega)} \qquad (3)$$

If convection dominated flow (or reaction dominated flow), then $M/\alpha \gg 1$:
$\rightarrow$ trade-off between $M/\alpha$ and $h^r$
$\rightarrow$ numerically prohibitive

$$\mathbb{P}e = h\frac{M}{\alpha}$$

Peclet number – "moral definition".
More precise definition an a
case-by-case mode.

$\rightarrow \mathbb{P}e$ should be less than 1 (for stability issues – see later).

**Idea**: stabilize the Galerkin method

*1D case*: Upwind method $\iff$ Artificial diffusion

*2D case*: Streamline diffusion (see lab)

$$+c(h) \int_\Omega \frac{1}{\|\mathbf{b}\|} \left(\mathbf{b} \cdot \nabla u_h\right) \left(\mathbf{b} \cdot \nabla v_h\right)$$

Artificial diffusion (diffuse everywhere)

$$+c(h) \int_\Omega \nabla u_h \cdot \nabla v_h$$

Not fully/strongly consistent!

$$\begin{cases} \text{Find } u_h \in V_h \\ a(u_h, v_h) + \mathscr{L}_h(u_h, f; v_h) = F(v_h) \qquad \forall \, v_h \in V_h \end{cases} \qquad (4)$$

$\mathscr{L}_h$ suitably chosen – should be such that:

$$\mathscr{L}_h(u, f; v_h) = 0 \qquad \forall \, v_h \in V_h$$

$\rightarrow$ **strongly consistent** approximation of the original problem.

**Idea**: proportional to the residual

$$\mathscr{L}_h(u_h, f; v_h) = \sum_{K \in \mathcal{T}_h} \int_K (\mathcal{L} u_h - f) \, \tau_K \phi(v_h) \qquad \forall \, v_h \in V_h \qquad (5)$$

$\tau_K$: scaling factor. Typical choice:

$$\tau_K(\mathbf{x}) = \delta \frac{h_K}{|\mathbf{b}(\mathbf{x})|} \qquad \forall \, \mathbf{x} \in K, K \in \mathcal{T}_h$$

$$h_K = \text{diam}(K) \qquad (6)$$

Many ways of choosing $\phi(v_h)$. Two remarkable choices:

1. $\phi(v_h) = \mathcal{L}v_h \rightarrow$ **GLS** – Galerkin least square method
2. $\phi(v_h) = \mathcal{L}_{\mathsf{ss}}v_h \rightarrow$ **SUPG** – Streamline upwind Petrov-Galerkin method

Notation: $\mathcal{L} = \mathcal{L}_{\mathsf{s}} + \mathcal{L}_{\mathsf{ss}}$ (symmetric + skew-symmetric part of $\mathcal{L}$)

Definitions:

$$_{V'}\langle \mathcal{L}_{\mathsf{s}}u, v\rangle_V = {}_V\langle u, \mathcal{L}_{\mathsf{s}}v\rangle_{V'} \qquad \forall\, u, v \in V$$
$$_{V'}\langle \mathcal{L}_{\mathsf{ss}}u, v\rangle_V = -{}_V\langle u, \mathcal{L}_{\mathsf{ss}}v\rangle_{V'} \qquad \forall\, u, v \in V$$

**Remark**: for matrices, $A = A_{\mathsf{S}} + A_{\mathsf{SS}}$, with:

$$A_{\mathsf{S}} = \frac{1}{2}(A + A^T), \qquad A_{\mathsf{SS}} = \frac{1}{2}(A - A^T)$$

## Example (Non-conservative form)

$$\mathcal{L}^1 = -\mu\Delta u + \mathbf{b}\cdot\nabla u + \sigma u$$

$$= \underbrace{\left[-\mu\Delta u + \left(\sigma - \frac{1}{2}\operatorname{div}\mathbf{b}\right)u\right]}_{\mathcal{L}_{\mathsf{s}}^1 u} + \underbrace{\left[\frac{1}{2}\left(\operatorname{div}(\mathbf{b}u) + \mathbf{b}\cdot\nabla u\right)\right]}_{\mathcal{L}_{\mathsf{ss}}^1 u}$$

Indeed:

$$_{V'}\langle \mathcal{L}_{\mathsf{s}}^1 u, v\rangle_V = \int_\Omega \mu\nabla u\cdot\nabla v + \left(\sigma - \frac{1}{2}\operatorname{div}\mathbf{b}\right)u\,v$$

$$= \int_\Omega \left[-\mu\Delta v + \left(\sigma - \frac{1}{2}\operatorname{div}\mathbf{b}\right)v\right]u = {}_V\langle u, \mathcal{L}_{\mathsf{s}}^1 v\rangle_{V'}$$

$$_{V'}\langle \mathcal{L}_{\mathsf{ss}}^1 u, v\rangle_V = \frac{1}{2}\int_\Omega \left(\operatorname{div}(\mathbf{b}u)v + (\mathbf{b}\cdot\nabla u)\,v\right)$$

$$= \frac{1}{2}\int_\Omega \left(-(\mathbf{b}u)\cdot\nabla v + (\mathbf{b}v)\cdot\nabla u\right)$$

$$= \frac{1}{2}\int_\Omega \left(-(\mathbf{b}\cdot\nabla v)u - \operatorname{div}(\mathbf{b}v)u\right) = -{}_V\langle u, \mathcal{L}_{\mathsf{ss}}^1 v\rangle_{V'}$$

## Example (Conservative form)

$$
\mathcal{L}^2 = -\mu\Delta u + \mathrm{div}(\mathbf{b}u) + \sigma u
$$

$$
= \underbrace{\left[-\mu\Delta u + \left(\sigma + \frac{1}{2}\,\mathrm{div}\,\mathbf{b}\right)u\right]}_{\mathcal{L}^2_s u} + \underbrace{\left[\frac{1}{2}\left(\mathrm{div}(\mathbf{b}u) + \mathbf{b}\cdot\nabla u\right)\right]}_{\mathcal{L}^2_{ss} u}
$$

The proof is similar (do it yourself).

## Remark

If $\operatorname{div} \mathbf{b} = 0$ (this happens, for instance, when $\mathbf{b}$ is constant), then the conservative and non-conservative forms coincide: $\mathcal{L}^1 = \mathcal{L}^2$. Indeed.

$$\operatorname{div}(\mathbf{b}u) = \mathbf{b} \cdot \nabla u$$

In this case:

$$\mathcal{L}_{\mathsf{s}}u = -\mu\Delta u + \sigma u, \qquad \mathcal{L}_{\mathsf{ss}}u = \mathbf{b} \cdot \nabla u$$

Indeed:

$$\begin{aligned}
{}_{V'}\langle \mathcal{L}_{\mathsf{s}}^1 u, v \rangle_V &= (\mu\nabla u, \nabla v) + (\sigma u, v) = {}_V\langle u, \mathcal{L}_{\mathsf{s}}^1 v \rangle_{V'} \\
{}_{V'}\langle \mathcal{L}_{\mathsf{ss}}^1 u, v \rangle_V &= (\mathbf{b} \cdot \nabla u, v) = (\nabla u, \mathbf{b}v) \\
&= -\left( u, \operatorname{div}(\mathbf{b}v) \right) = -\left( u, \mathbf{b} \cdot \nabla v \right) = -{}_V\langle u, \mathcal{L}_{\mathsf{ss}}^1 v \rangle_{V'}
\end{aligned}$$

Back to stabilized Galerkin.

### Remark

Note that if $r = 1$, $\sigma = 0$ and div $\mathbf{b} = 0$, the two methods SUPG and GLS coincide. Indeed, $-\Delta u_h|_K = 0$ on each $K \in \mathcal{T}_h$.

**Problem in conservative form – GLS method**

$$
\begin{cases}
\text{Find } u_h \in V_h \\
a(u_h, v_h) + \displaystyle\sum_{K \in \mathcal{T}_h} \int_K \mathcal{L} u_h \, \tau_K \, \mathcal{L} v_h = \\
\displaystyle\int_\Omega f \, v_h + \sum_{K \in \mathcal{T}_h} \int_K f \, \tau_K \, \mathcal{L} v_h \qquad \forall \, v_h \in V_h
\end{cases}
\tag{7}
$$

which can be rewritten (with obvious meaning of notations) as:

$$
\begin{cases}
\text{Find } u_h \in V_h \\
a_h(u_h, v_h) = F_h(v_h) \qquad \forall \, v_h \in V_h
\end{cases}
\tag{8}
$$

# Stability analysis

## Theorem

*Consider the conservative case. Suppose that*

$$\exists \, \gamma_0, \gamma_1 > 0 : 0 < \gamma_0 \leq \gamma(\mathbf{x}) \leq \gamma_1 \qquad (9)$$

*then, for a suitable constant C independent of h, we have:*

$$\|u_h\|_{GLS}^2 \leq C \, \|f\|_{L^2(\Omega)}^2$$

*($\|\cdot\|_{GLS}$ to be defined later).*

**Proof.** Take $v_h = u_h$. We have:

$$a_h(u_h, u_h) = \int_\Omega \mu |\nabla u_h|^2 + \underbrace{\int_\Omega \text{div}(\mathbf{b} u_h) u_h}_{=-\int_\Omega \mathbf{b}\cdot(u_h\nabla u_h)=-\frac{1}{2}\int_\Omega \mathbf{b}\cdot\nabla(u_h^2)=\frac{1}{2}\int_\Omega \text{div}(\mathbf{b})u_h^2} + \int_\Omega \sigma u_h^2 + \sum_{K\in\mathcal{T}_h}\int_K \tau_K \left(\mathcal{L}u_h\right)^2$$

$$= \int_\Omega \mu |\nabla u_h|^2 + \int_\Omega \underbrace{\left(\sigma + \frac{1}{2}\text{div}\,\mathbf{b}\right)}_{=:\gamma(\mathbf{x})} u_h^2 + \sum_{K\in\mathcal{T}_h}\int_K \tau_K \left(\mathcal{L}u_h\right)^2$$

$$=: \|u_h\|^2_{\text{GLS}}$$

On the other hand:

$$|F_h(u_h)| \leq \left| \int_\Omega f \, u_h \right| + \left| \sum_{K \in \mathcal{T}_h} \int_K f \, \tau_K \, \mathcal{L} u_h \right|$$

Where:

$$
\begin{aligned}
\left| \int_\Omega f \, u_h \right| &= \left| \int_\Omega \frac{1}{\sqrt{\gamma}} f \, \sqrt{\gamma} u_h \right| \\
&\overset{\text{(Cauchy–Schwarz)}}{\leq} \left\| \frac{1}{\sqrt{\gamma}} f \right\|_{L^2(\Omega)} \| \sqrt{\gamma} u_h \|_{L^2(\Omega)} \\
&\overset{\text{(Young*)}}{\leq} \left\| \frac{1}{\sqrt{\gamma}} f \right\|_{L^2(\Omega)}^2 + \frac{1}{4} \| \sqrt{\gamma} u_h \|_{L^2(\Omega)}^2
\end{aligned}
$$

---

* Young inequality: $AB \leq \epsilon A^2 + \frac{1}{4\epsilon} B^2 \quad \forall A, B \in \mathbb{R}, \epsilon > 0$

And where:

$$\left| \sum_{K \in \mathcal{T}_h} \int_K f \, \tau_K \, \mathcal{L} u_h \right| = \left| \sum_{K \in \mathcal{T}_h} \int_K \sqrt{\tau_K} f \, \sqrt{\tau_K} \mathcal{L} u_h \right|$$

$$\overset{\text{(Cauchy–Schwarz)}}{\leq} \sum_{K \in \mathcal{T}_h} \| \sqrt{\tau_K} f \|_{L^2(K)} \, \| \sqrt{\tau_K} \mathcal{L} u_h \|_{L^2(K)}$$

$$\overset{\text{(Young)}}{\leq} \sum_{K \in \mathcal{T}_h} \| \sqrt{\tau_K} f \|_{L^2(K)}^2 + \frac{1}{4} \| \sqrt{\tau_K} \mathcal{L} u_h \|_{L^2(K)}^2$$

To wrap-up, $a_h(u_h, u_h) = F_h(u_h)$ implies:

$$\|u_h\|_{\text{GLS}}^2 = \int_\Omega \mu|\nabla u_h|^2 + \int_\Omega \gamma\, u_h^2 + \sum_{K \in \mathcal{T}_h} \int_K \tau_K \left(\mathcal{L}u_h\right)^2$$

$$\leq \left[ \left\| \frac{1}{\sqrt{\gamma}} f \right\|_{L^2(\Omega)}^2 + \sum_{K \in \mathcal{T}_h} \|\sqrt{\tau_K} f\|_{L^2(\Omega)}^2 \right]$$

$$+ \frac{1}{4} \left[ \int_\Omega \gamma\, u_h^2 + \sum_{K \in \mathcal{T}_h} \int_\Omega \tau_K \left(\mathcal{L}u_h\right)^2 \right]$$

$$\leq \underbrace{\left( \frac{1}{\gamma_0} + \max_{K \in \mathcal{T}_h} \tau_K \right)}_{C \text{ (if } \tau_K \text{ uniformly bounded w.r.t. } h)} \|f\|_{L^2(\Omega)}^2 + \frac{1}{4}\|u_h\|_{\text{GLS}}^2$$

$\rightarrow$ $\qquad \boxed{\|u_h\|_{\text{GLS}}^2 \leq \frac{4}{3} C\, \|f\|_{L^2(\Omega)}^2}$ **Stability**

# On the choice of $\tau_K$ (stabilization parameter)

First choice is (6): $\tau_K(\mathbf{x}) = \delta \dfrac{h_K}{|\mathbf{b}(\mathbf{x})|}$, with $\delta > 0$ to be chosen.

Alternative choice:

$$\tau_K(\mathbf{x}) = \frac{h_K}{2|\mathbf{b}(\mathbf{x})|} \xi(\mathbb{P}e_K) \tag{10}$$

where:

$\xi(\theta) = \coth(\theta) - \dfrac{1}{\theta}$

$\mathbb{P}e_K(\mathbf{x}) = \dfrac{|\mathbf{b}(\mathbf{x})|}{2\,\mu(\mathbf{x})} h_K$
(local Peclet number)

were $\dfrac{|\mathbf{b}(\mathbf{x})|}{2\,\mu(\mathbf{x})}$ is the global/physical Peclet number

### Remark

Note that since $\lim_{\theta \to +\infty} \xi(\theta) = 1$, if $\mathbb{P}e_K(\mathbf{x}) \gg 1$, $\tau_K$ defined in (10) reduces, in the limit, to (6) with $\delta = 1/2$.

Moreover, if $\theta \to 0$, then $\xi(\theta) = \theta/3 + o(\theta)$, therefore when $\mathbb{P}e_K(\mathbf{x}) \ll 1$, we have $\tau_K(\mathbf{x}) \to 0$ and **no stabilization** is needed (pure Galerkin works fine!).

### Remark

Note also (more tricky) that, in 1D, choice (10) coincides with the famous **Scharfetter-Gummel** stabilization scheme, a second order scheme that is nodally exact (see [Q], Secs. 13.8.7 and 13.6).

---

[Q] A. Quarteroni, *Numerical Models for Differential Problems*, 3rd Ed., Springer, 2018

## Remark: polynomial of degree $r \geq 1$ arbitrary

(10) modifies as:

$$\mathbb{P}e_K^r = \frac{|\mathbf{b}(\mathbf{x})|}{2\,\mu(\mathbf{x})\,r} h_K$$

and then:

$$\tau_K(\mathbf{x}) = \frac{h_K}{2\,|\mathbf{b}(\mathbf{x})|\,r} \xi(\mathbb{P}e_K^r(\mathbf{x}))$$

# Convergence of GLS

To state the convergence result for GLS, we need the following inequality, known as *inverse inequality*.

> **Inverse inequality**
>
> $$\sum_{K \in \mathcal{T}_h} h_K^2 \int_K |\Delta v_h|^2 dK \leq C_0 \|\nabla v_h\|_{L^2(\Omega)}^2 \quad \forall v_h \in X_h^r. \tag{11}$$

### Theorem (Convergence of GLS)

*Assume that the space $V_h$ satisfies the following local approximation property: for each $v \in V \cap H^{r+1}(\Omega)$, there exists a function $\hat{v}_h \in V_h$ s.t.*

$$\|v - \hat{v}_h\|_{\mathrm{L}^2(K)} + h_K |v - \hat{v}_h|_{\mathrm{H}^1(K)} + h_K^2 |v - \hat{v}_h|_{\mathrm{H}^2(K)} \leq C h_K^{r+1} |v|_{\mathrm{H}^{r+1}(K)} \quad (12)$$

*for each $K \in \mathcal{T}_h$. Moreover, we suppose that for each $K \in \mathcal{T}_h$ the local Péclet number of $K$ satisfies*

$$\mathbb{P}e_K(\mathbf{x}) = \frac{|\mathbf{b}(\mathbf{x})| \, h_K}{2\mu} > 1 \quad \forall \mathbf{x} \in K, \tag{13}$$

*that is, we are in the pre-asymptotic regime. Finally, we suppose that the inverse inequality holds and that the stabilization parameter satisfies the relation $0 < \delta \leq 2C_0^{-1}$.*
*Then, as long as $u \in H^{r+1}(\Omega)$, the following super-optimal estimate holds:*

$$\|u_h - u\|_{GLS} \leq C h^{r+1/2} |u|_{\mathrm{H}^{r+1}(\Omega)}. \tag{14}$$

**Proof.**
First of all, we rewrite the error as follows

$$e_h = u_h - u = \sigma_h - \eta, \tag{15}$$

with $\sigma_h = u_h - \hat{u}_h$, $\eta = u - \hat{u}_h$, where $\hat{u}_h \in V_h$ is a function that depends on $u$ and that satisfies property (12). If, for instance, $V_h = X_h^r \cap H_0^1(\Omega)$, we can choose $\hat{u}_h = \Pi_h^r u$, that is the finite element interpolant of $u$.

We start by estimating the norm $\|\sigma_h\|_{GLS}$. By exploiting the strong consistency of the GLS scheme, we obtain

$$||\sigma_h||_{GLS}^2 = a_h(\sigma_h, \sigma_h) = a_h(u_h - u + \eta, \sigma_h) = a_h(\eta, \sigma_h).$$

Now, thanks to the homogeneous Dirichlet boundary conditions it follows that, by adding and subtracting the term $\sum_{K \in \mathcal{T}_h} (\eta, \mathcal{L}\sigma_h)_K$, suitable computations lead to:

$$a_h(\eta, \sigma_h) = \mu \int_\Omega \nabla \eta \cdot \nabla \sigma_h \, d\Omega - \int_\Omega \eta \, \mathbf{b} \cdot \nabla \sigma_h \, d\Omega + \int_\Omega \sigma \, \eta \, \sigma_h \, d\Omega$$

$$+ \sum_{K \in \mathcal{T}_h} \delta \Big( \mathcal{L}\eta, \frac{h_K}{|\mathbf{b}|} \mathcal{L}\sigma_h \Big)_{L^2(K)}$$

$$= \underbrace{\mu(\nabla \eta, \nabla \sigma_h)_{L^2(\Omega)}}_{\text{(I)}} - \underbrace{\sum_{K \in \mathcal{T}_h} (\eta, \mathcal{L}\sigma_h)_{L^2(K)}}_{\text{(II)}} + \underbrace{2(\gamma \, \eta, \sigma_h)_{L^2(\Omega)}}_{\text{(III)}}$$

$$+ \underbrace{\sum_{K \in \mathcal{T}_h} (\eta, -\mu \Delta \sigma_h)_{L^2(K)}}_{\text{(IV)}} + \underbrace{\sum_{K \in \mathcal{T}_h} \delta \Big( \mathcal{L}\eta, \frac{h_K}{|\mathbf{b}|} \mathcal{L}\sigma_h \Big)_{L^2(K)}}_{\text{(V)}}.$$

We now bound the terms (I)-(V) separately.

By carefully using the Cauchy-Schwarz and Young inequalities we obtain

$$|(\text{I})| = |\mu(\nabla\eta, \nabla\sigma_h)_{\mathrm{L}^2(\Omega)}| \leq \frac{\mu}{4}\|\nabla\sigma_h\|^2_{\mathrm{L}^2(\Omega)} + \mu\|\nabla\eta\|^2_{\mathrm{L}^2(\Omega)},$$

$$\begin{aligned}
|(\text{II})| &= \left|\sum_{K\in\mathcal{T}_h}(\eta, L\sigma_h)_{\mathrm{L}^2(K)}\right| \\
&= \left|\sum_{K\in\mathcal{T}_h}\left(\sqrt{\frac{|\mathbf{b}|}{\delta\,h_K}}\,\eta, \sqrt{\frac{\delta\,h_K}{|\mathbf{b}|}}\,L\sigma_h\right)_{\mathrm{L}^2(K)}\right| \\
&\leq \frac{1}{4}\sum_{K\in\mathcal{T}_h}\delta\left(\frac{h_K}{|\mathbf{b}|}L\sigma_h, L\sigma_h\right)_{\mathrm{L}^2(K)} + \sum_{K\in\mathcal{T}_h}\left(\frac{|\mathbf{b}|}{\delta\,h_K}\eta, \eta\right)_{\mathrm{L}^2(K)},
\end{aligned}$$

$$\begin{aligned}
|(\text{III})| &= 2|(\gamma\,\eta, \sigma_h)_{\mathrm{L}^2(\Omega)}| = 2|(\sqrt{\gamma}\,\eta, \sqrt{\gamma}\,\sigma_h)_{\mathrm{L}^2(\Omega)}| \\
&\leq \frac{1}{2}\|\sqrt{\gamma}\,\sigma_h\|^2_{\mathrm{L}^2(\Omega)} + 2\|\sqrt{\gamma}\,\eta\|^2_{\mathrm{L}^2(\Omega)}.
\end{aligned}$$

For the term $(\mathrm{IV})$, thanks again to the Cauchy-Schwarz and Young inequalities, hypothesis (13) and the inverse inequality (11), we obtain

$$
\begin{aligned}
|(\mathrm{IV})| &= \left| \sum_{K \in \mathcal{T}_h} (\eta, -\mu \Delta \sigma_h)_{\mathrm{L}^2(K)} \right| \\
&\leq \frac{1}{4} \sum_{K \in \mathcal{T}_h} \delta \, \mu^2 \Big( \frac{h_K}{|\mathbf{b}|} \, \Delta \sigma_h, \Delta \sigma_h \Big)_{\mathrm{L}^2(K)} \\
&\quad + \sum_{K \in \mathcal{T}_h} \Big( \frac{|\mathbf{b}|}{\delta \, h_K} \eta, \eta \Big)_{\mathrm{L}^2(K)} \\
&\leq \frac{1}{8} \, \delta \, \mu \sum_{K \in \mathcal{T}_h} h_K^2 \, (\Delta \sigma_h, \Delta \sigma_h)_{\mathrm{L}^2(K)} + \sum_{K \in \mathcal{T}_h} \Big( \frac{|\mathbf{b}|}{\delta \, h_K} \eta, \eta \Big)_{\mathrm{L}^2(K)} \\
&\leq \frac{\delta \, C_0 \, \mu}{8} \| \nabla \sigma_h \|_{\mathrm{L}^2(\Omega)}^2 + \sum_{K \in \mathcal{T}_h} \Big( \frac{|\mathbf{b}|}{\delta \, h_K} \eta, \eta \Big)_{\mathrm{L}^2(K)}.
\end{aligned}
$$

Term (V) can finally be bounded once again thanks to the Cauchy-Schwarz and Young inequalities as follows

$$
\begin{aligned}
|(\mathrm{V})| &= \left| \sum_{K \in \mathcal{T}_h} \delta \left( L\eta, \frac{h_K}{|\mathbf{b}|} L\sigma_h \right)_{\mathrm{L}^2(K)} \right| \\
&\leq \frac{1}{4} \sum_{K \in \mathcal{T}_h} \delta \left( \frac{h_K}{|\mathbf{b}|} L\sigma_h, L\sigma_h \right)_{\mathrm{L}^2(K)} + \sum_{K \in \mathcal{T}_h} \delta \left( \frac{h_K}{|\mathbf{b}|} L\eta, L\eta \right)_{\mathrm{L}^2(K)}.
\end{aligned}
$$

Thanks to these upper bounds, we obtain the following estimate

$$\|\sigma_h\|_{GLS}^2 = a_h(\eta, \sigma_h) \leq \frac{1}{4} \|\sigma_h\|_{GLS}^2$$

$$+ \frac{1}{4} \left( \|\sqrt{\gamma}\, \sigma_h\|_{L^2(\Omega)}^2 + \sum_{K \in \mathcal{T}_h} \delta \Big( \frac{h_K}{|\mathbf{b}|}\, L\sigma_h, L\sigma_h \Big)_{L^2(K)} \right) + \frac{\delta\, C_0\, \mu}{8} \|\nabla \sigma_h\|_{L^2(\Omega)}^2$$

$$+ \underbrace{\mu\, \|\nabla \eta\|_{L^2(\Omega)}^2 + 2 \sum_{K \in \mathcal{T}_h} \Big( \frac{|\mathbf{b}|}{\delta\, h_K}\eta, \eta \Big)_{L^2(K)} + 2\, \|\sqrt{\gamma}\, \eta\|_{L^2(\Omega)}^2 + \sum_{K \in \mathcal{T}_h} \delta \Big( \frac{h_K}{|\mathbf{b}|}\, L\eta, L\eta \Big)_{L^2(K)}}_{\mathcal{E}(\eta)}$$

$$\leq \frac{1}{2} \|\sigma_h\|_{GLS}^2 + \mathcal{E}(\eta),$$

having exploited, in the last passage, the assumption that $\delta \leq 2 C_0^{-1}$.

Then, we can state that

$$\|\sigma_h\|_{GLS}^2 \leq 2\, \mathcal{E}(\eta).$$

We now estimate the term $\mathcal{E}(\eta)$, by bounding each of its summands separately. To this end, we will basically use the local approximation property (12) and the requirement formulated in (13) on the local Péclet number $\mathbb{P}e_K$.

Moreover, we observe that the constants $C$, introduced in the remainder, depend neither on $h$ nor on $\mathbb{P}e_K$, but can depend on other quantities such as the constant $\gamma_1$ in (9), the reaction constant $\sigma$, the norm $||\mathbf{b}||_{L^\infty(\Omega)}$, the stabilization parameter $\delta$.

We then have

$$
\begin{aligned}
\mu \, \|\nabla\eta\|_{\mathrm{L}^2(\Omega)}^2 &\le C \, \mu \, h^{2r} \, |u|_{\mathrm{H}^{r+1}(\Omega)}^2 \\
&\le C \, \frac{||\mathbf{b}||_{L^\infty(\Omega)} \, h}{2} \, h^{2r} \, |u|_{\mathrm{H}^{r+1}(\Omega)}^2 \le C \, h^{2r+1} \, |u|_{\mathrm{H}^{r+1}(\Omega)}^2
\end{aligned}
$$

$$2 \sum_{K \in \mathcal{T}_h} \left( \frac{|\mathbf{b}|}{\delta \, h_K} \eta, \eta \right)_{\mathrm{L}^2(K)} \leq C \, \frac{\|\mathbf{b}\|_{L^\infty(\Omega)}}{\delta} \sum_{K \in \mathcal{T}_h} \frac{1}{h_K} \, h_K^{2(r+1)} \, |u|_{\mathrm{H}^{r+1}(K)}^2$$
$$\leq C \, h^{2r+1} \, |u|_{\mathrm{H}^{r+1}(\Omega)}^2,$$

$$2 \, \|\sqrt{\gamma} \, \eta\|_{\mathrm{L}^2(\Omega)}^2 \leq 2 \, \gamma_1 \, \|\eta\|_{\mathrm{L}^2(\Omega)}^2 \leq C \, h^{2(r+1)} \, |u|_{\mathrm{H}^{r+1}(\Omega)}^2,$$

having exploited, for controlling the third summand, the assumption (9).

Finding an upper bound for the fourth summand of $\mathcal{E}(\eta)$ is slightly more difficult: first, by elaborating on the term $L\eta$, we have

$$
\sum_{K \in \mathcal{T}_h} \delta \left( \frac{h_K}{|\mathbf{b}|} L\eta, L\eta \right)_{\mathrm{L}^2(K)}
$$

$$
= \sum_{K \in \mathcal{T}_h} \delta \left\| \sqrt{\frac{h_K}{|\mathbf{b}|}} L\eta \right\|_{\mathrm{L}^2(K)}^2
$$

$$
= \sum_{K \in \mathcal{T}_h} \delta \left\| - \mu \sqrt{\frac{h_K}{|\mathbf{b}|}} \Delta\eta + \sqrt{\frac{h_K}{|\mathbf{b}|}} \operatorname{div}(\mathbf{b}\eta) + \sigma \sqrt{\frac{h_K}{|\mathbf{b}|}} \eta \right\|_{\mathrm{L}^2(K)}^2 \qquad (16)
$$

$$
\leq C \sum_{K \in \mathcal{T}_h} \delta \left( \left\| \mu \sqrt{\frac{h_K}{|\mathbf{b}|}} \Delta\eta \right\|_{\mathrm{L}^2(K)}^2 + \left\| \sqrt{\frac{h_K}{|\mathbf{b}|}} \operatorname{div}(\mathbf{b}\eta) \right\|_{\mathrm{L}^2(K)}^2 \right.
$$

$$
\left. + \left\| \sigma \sqrt{\frac{h_K}{|\mathbf{b}|}} \eta \right\|_{\mathrm{L}^2(K)}^2 \right).
$$

Now, with a similar computation to the one performed to obtain estimates (28) and (29), it is easy to prove that the second and third summands of the left-hand side of (16) can be bounded using a term of the form $C\,h^{2r+1}\,|u|^2_{\mathrm{H}^{r+1}(\Omega)}$, for a suitable choice of the constant $C$. For the first summand, we have

$$
\sum_{K\in\mathcal{T}_h} \delta \left\| \mu\,\sqrt{\frac{h_K}{|\mathbf{b}|}}\,\Delta\eta \right\|^2_{\mathrm{L}^2(K)} \leq \sum_{K\in\mathcal{T}_h} \delta\,\frac{h_K^2\,\mu}{2}\,\|\Delta\eta\|^2_{\mathrm{L}^2(K)}
$$
$$
\leq C\,\delta\,\|\mathbf{b}\|_{L^\infty(\Omega)} \sum_{K\in\mathcal{T}_h} h_K^3\,\|\Delta\eta\|^2_{\mathrm{L}^2(K)} \leq C\,h^{2r+1}\,|u|^2_{\mathrm{H}^{r+1}(\Omega)},
$$

having again used conditions (12) and (13). The latter bound allows us to conclude that

$$
\mathcal{E}(\eta) \leq C\,h^{2r+1}\,|u|^2_{\mathrm{H}^{r+1}(\Omega)},
$$

that is

$$
\|\sigma_h\|_{GLS} \leq C\,h^{r+1/2}\,|u|_{\mathrm{H}^{r+1}(\Omega)}. \tag{17}
$$

Reverting to (15), to obtain the desired estimate for the norm $\|u_h - u\|_{GLS}$ we still have to estimate $\|\eta\|_{GLS}$. This evidently leads to estimating three contributions as in (28), (29) and (16), and eventually produces

$$\|\eta\|_{GLS} \leq C \, h^{r+1/2} \, |u|_{\mathrm{H}^{r+1}(\Omega)}.$$

The desired estimate (14) follows by combining this result with (17). $\quad\square$

Rationale for (13): use GLS only beyond the asymptotic regime.
In the asymptotic regime, i.e. if $\mathbb{P}e_K < 1$, use standard Galerkin.

Notice that (14) is **"super-optimal"**: $h^{r+1/2}$ instead of $h^r$ as in Galerkin.
The reason is (13) (a lower limit for $h$).

We are not in the asymptotic regime for $h \to 0$ ( $\implies \mathbb{P}e_K \to 0$).



$$\mathcal{O}(h^r) \qquad\qquad \mathcal{O}(h^{r+1/2})$$

$\mathbb{P}e_K(\mathbf{x})$

1

Galerkin optimal     Stabilized Galerkin (GLS
                     or SUPG). Galerkin
                     would be unstable

⇓                         ⇓

use Galerkin         use stabilized Galerkin

[Q] A. Quarteroni, *Numerical Models for Differential Problems*, 3rd Ed., Springer, 2018

The Discontinuous Galerkin method can be extended to the diffusion-transport-reaction problem in conservation form:

$$\begin{cases} -\operatorname{div}(\mu \nabla u + \mathbf{b}u) + \sigma u = f & \text{in } \Omega \\ u = 0 & \text{on } \partial\Omega \end{cases} \tag{18}$$

We introduce the following space:

$$W_\delta^0 = \left\{ v_\delta \in L^2(\Omega) \colon \ v_\delta|_K \in H^1(K) \, \forall \, K \in \mathcal{T}_h, \ v_\delta|_{\partial\Omega} = 0 \right\}$$

The Discontinuous Galerkin formulation reads as follow:
find $u_\delta \in W_\delta^0$ s.t.

$$
\sum_{K \in \mathcal{T}_h} (\mu \nabla u_\delta, \nabla v_\delta)_{L^2(K)} - \sum_{e \in \mathcal{E}_\delta} \int_e [\![v_\delta]\!] \cdot \{\!\{\mu \nabla u_\delta\}\!\} - \theta \sum_{e \in \mathcal{E}_\delta} \int_e [\![u_\delta]\!] \{\!\{\mu \nabla v_\delta\}\!\}
$$
$$
+ \sum_{e \in \mathcal{E}_\delta} \int_e \overline{\gamma} \, [\![u_\delta]\!] \cdot [\![v_\delta]\!] - \sum_{K \in \mathcal{T}_h} (\mathbf{b} u_\delta, \nabla v_\delta)_{L^2(K)} + \sum_{e \in \mathcal{E}_\delta} \int_e \{\!\{\mathbf{b} u_\delta\}\!\}_{\mathbf{b}} \cdot [\![v_\delta]\!]
$$
$$
+ \sum_{K \in \mathcal{T}_h} (\sigma u_\delta, v_\delta)_{L^2(K)} = \sum_{K \in \mathcal{T}_h} (f, v_\delta)_{L^2(K)} \ ,
$$

$$(19)$$

where $\mathcal{E}_\delta \equiv \mathcal{F}_h$ is the set of the edges of the elements $\{K\}$, $\overline{\gamma}$ is the DG stabilization function (see Lecture 03, slide 12) and where

$$
\{\!\{\mathbf{b} u_\delta\}\!\}_{\mathbf{b}} = \begin{cases} \mathbf{b} u_\delta^+ & \text{if } \mathbf{b} \cdot \mathbf{n}^+ > 0 \\ \mathbf{b} u_\delta^- & \text{if } \mathbf{b} \cdot \mathbf{n}^+ < 0 \\ \mathbf{b} \{\!\{u_\delta\}\!\} & \text{if } \mathbf{b} \cdot \mathbf{n}^+ = 0 \ . \end{cases}
$$

$$(20)$$

Observe that $\{\!\{\mathbf{b} u_\delta\}\!\}_{\mathbf{b}} \cdot [\![v_\delta]\!] = 0$ if $\mathbf{b} \cdot \mathbf{n}^+ = 0$.

If the diffusion-transport-reaction equation is written in non-conservative form, by $\mathbf{b} \cdot \nabla u = \operatorname{div}(\mathbf{b}u) - \operatorname{div}(\mathbf{b})u$ it is sufficient to modify (19) by substituting the term

$$\sum_{m=1}^{M} (\sigma u_\delta, v_\delta)_{\Omega_m} \quad \text{with} \quad \sum_{m=1}^{M} (\eta u_\delta, v_\delta)_{\Omega_m},$$

where $\eta(\mathbf{x}) = \sigma(\mathbf{x}) - \operatorname{div}(\mathbf{b}(\mathbf{x}))$.

This time we suppose that there exists a positive constant $\eta_0 > 0$ so that $\eta(\mathbf{x}) \geq \eta_0$ for almost every $\mathbf{x} \in \Omega$.

We now present some numerical solutions obtained using linear finite elements for the following two-dimensional diffusion-transport problem

$$\begin{cases} -\mu \Delta u + \mathbf{b} \cdot \nabla u = f & \text{in } \Omega = (0,1) \times (0,1), \\ u = g & \text{on } \partial\Omega, \end{cases} \tag{21}$$

where $\mathbf{b} = (-1, 1)^T$. To start with let us consider the following constant data: $f \equiv 1$ and $g \equiv 0$. In this case the solution is characterized by a boundary layer near the edges $x = 0$ and $y = 1$.

Figure: Approximation of problem (21) with $\mu = 10^{-3}$, $h = 1/80$, using the standard (left) and GLS (right) Galerkin method. The corresponding local Péclet number is $\mathbb{P}e_K = 8.84$

Figure: Approximation of problem (21) with $\mu = 10^{-3}$, $h = 1/20$, using the standard (left) and GLS (right) Galerkin method. The corresponding local Péclet number is $\mathbb{P}e_K = 35.35$

Figure: Approximation of problem (21) with $\mu = 10^{-5}$, $h = 1/80$, using the standard (left) and GLS (right) Galerkin method. The corresponding local Péclet number is $\mathbb{P}e_K = 883.88$

Figure: Approximation of problem (21) with $\mu = 10^{-5}$, $h = 1/20$, using the standard (left) and GLS (right) Galerkin method. The corresponding local Péclet number is $\mathbb{P}e_K = 3535.5$

Let us now set $\mathbf{b} = (1,1)^T$ and choose forcing term $f$ and the boundary data $g$ in such a way that

$$u(x,y) = x + y(1-x) + \frac{e^{-1/\mu} - e^{-(1-x)(1-y)/\mu}}{1 - e^{-1/\mu}}$$

is the exact solution.

The corresponding Péclet number is $\mathbb{P}e = (\sqrt{2}\mu)^{-1}$.

For small values of the viscosity $\mu$, this solution features a boundary layer near the edges $x = 1$ and $y = 1$.

(a) Standard Galerkin method    (b) SUPG method



(c) DG method

Figure: The approximate solution of problem (21) for $\mu = 10^{-1}$. Triangular grid with $h \approx 1/8$ and piecewise linear finite elements ($r = 1$).
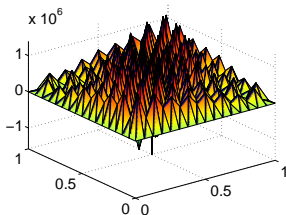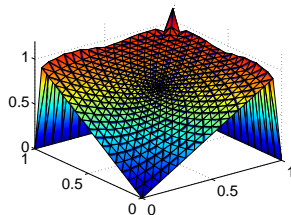
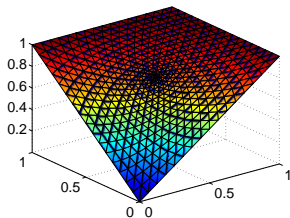(a) Standard Galerkin method

(b) SUPG method



(c) DG method

Figure: The approximate solution of problem (21) for $\mu = 10^{-9}$. Triangular grid with $h \approx 1/8$ and piecewise linear finite elements ($r = 1$).

(a) Standard Galerkin method



(b) SUPG method



(c) DG method

Figure: The approximate solution of problem (21) for $\mu = 10^{-1}$. Triangular grid with $h \approx 1/16$ and piecewise linear finite elements ($r = 1$).

(a) Standard Galerkin method

(b) SUPG method



(c) DG method

Figure: The approximate solution of problem (21) for $\mu = 10^{-9}$. Triangular grid with $h \approx 1/16$ and piecewise linear finite elements ($r = 1$).

Finally, we consider a pure transport problem, that is $\mathbf{b} \cdot \nabla u = f$ in $\Omega = (0,1)^2$ with $u = g$ su $\Gamma^-$, $\mathbf{b} = (1,1)$, with $f$ and $g$ chosen in such a way that the exact solution is $u(x,y) = 1 + \sin(\pi(x+1)(y+1)^2/8)$.
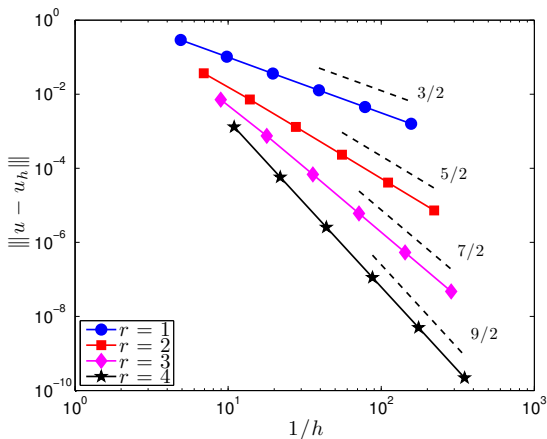
We solve this problem by the DG method with piecewise polynomials of degree $r = 1, 2, 3$ and $4$ on a sequence of uniform triangular grids with gridsize $h$.

The DG method provides the following error estimate [1]

$$
\begin{aligned}
\|u - u_h\| &= \left( \|u - u_h\|_{L^2(\Omega)}^2 + \sum_{e \in \mathcal{E}_h} \|s_e^{1/2}[u - u_h]\|_{0,e}^2 \right)^{1/2} \\
&\leq C h^{r+1/2} \|u\|_{H^{r+1}(\Omega)},
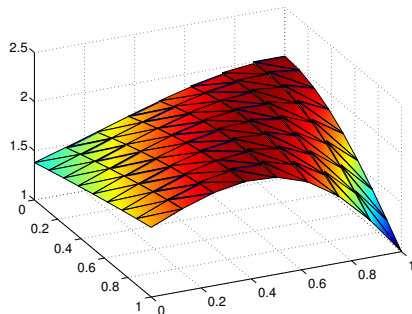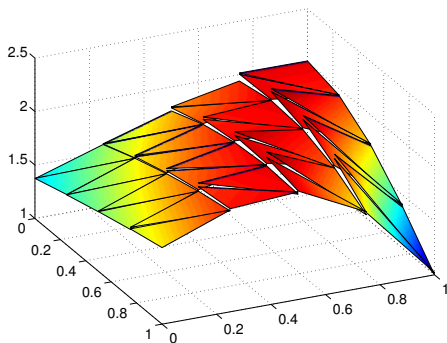\end{aligned}
\tag{22}
$$

where $s_e = \alpha |\mathbf{b} \cdot \mathbf{n_e}|$ is a suitable stabilization term, where $\alpha$ is a positive constant independent of $h$ and $e$. $\mathcal{E}_h$ is the set of all the edges of the triangulation and $C$ is a positive constant.

---

[1] see e.g. F. Brezzi, L. D. Marini and E. Süli, "Discontinuos Galerkin methods for first-order hyperbolic problems", *Math. Models Methods Appl. Sci.* (2004)

Figure: Approximation error (in the energy norm (22)) vs number of degrees of freedom for finite elements of degree $r = 1, 2, 3, 4$

Figure: Finite element solutions obtained on a uniform grid with gridsize $h = 1/4$ (left) and $h = 1/8$ (right)