

Chapter 2

Finite Element Methods (FEM)

2.1 Generalities

Remark 2.1. Finite element methods. Finite element methods were one of the main topics of Numerical Mathematics 3. The knowledge of the lecture notes of Numerical Mathematics 3 is assumed. Only a few issues, which are important for this course on numerical methods for convection-dominated problems, will be reminded here.

Let $\{\mathcal{T}^h\}$ be a family of regular triangulations consisting of mesh cells $\{K\}$. The triangulations are assumed to be quasi-uniform. The diameter of a mesh cell K is denoted by h_K and it is $h = \max_K \{h_K\}$. Parametric finite element spaces will be considered with affine maps between a reference cell \hat{K} and all physical cells K . \square

Theorem 2.2. Local interpolation error estimate. *Let $I_K : C^s(K) \rightarrow P(K)$ be an interpolation operator as defined in Numerical Mathematics 3, where $P(K)$ is a polynomial space defined on K . Let $p \in [1, \infty)$ and $(m+1-s)p > 1$. Then there is a constant C , which is independent of $v \in W^{m+1,p}(K)$, such that*

$$\|D^k(v - I_K v)\|_{L^p(K)} \leq C h_K^{m+1-k} \|D^{m+1}v\|_{L^p(K)}, \quad 0 \leq k \leq m+1. \quad (2.1)$$

for all $v \in W^{m+1,p}(K)$.

Proof. See lecture notes of Numerical Mathematics 3. \blacksquare

Theorem 2.3. Inverse estimate. *Let $0 \leq k \leq l$ be natural numbers and let $p, q \in [1, \infty]$. Then there is a constant C_{inv} , which depends only on $k, l, p, q, \hat{K}, \hat{P}(\hat{K})$ such that*

$$\|D^l v^h\|_{L^q(K)} \leq C_{\text{inv}} h_K^{(k-l)-d(p^{-1}-q^{-1})} \|D^k v^h\|_{L^p(K)} \quad \forall v^h \in P(K). \quad (2.2)$$

Proof. See lecture notes of Numerical Mathematics 3. \blacksquare

2.2 The Galerkin Method

Remark 2.4. On the size of the constant in the Lemma of Cea. The properties of the bilinear form from problem (1.27) were studied in the proof of Theorem 1.11. It was shown that with appropriate regularity assumptions and under condition (1.30), the bilinear form is bounded with a constant M of order $\max\{\|\mathbf{b}\|_{L^\infty(\Omega)}, \|\sigma\|_{L^\infty(\Omega)}\}$ and it is coercive a constant of order $m = \varepsilon$. In this case, one can apply the Lemma of Cea and one obtains the error estimate

$$\|u - u^h\|_V \leq \frac{C \max\{\|\mathbf{b}\|_{L^\infty(\Omega)}, \|\sigma\|_{L^\infty(\Omega)}\}}{\varepsilon} \inf_{v^h \in V^h} \|u - v^h\|_V, \quad C \in \mathbb{R}.$$

In the convection-dominated case $\varepsilon \ll L \|\mathbf{b}\|_{L^\infty(\Omega)}$, where L is a characteristic length scale of the problem, the first factor of this estimate becomes very large.

Thus, from this error estimate one cannot expect that the Galerkin finite element solution is accurate unless the second factor, which is the best approximation error, is very small. On uniformly refined grids, the best approximation error becomes very small only if the dimension of the finite element space V^h becomes very large. \square

Example 2.5 (Galerkin method). A standard test problem in two dimensions has the form

$$\begin{aligned} -\varepsilon \Delta u + (1, 0)^T \cdot \nabla u &= 1 \quad \text{in } \Omega = (0, 1)^2, \\ u &= 0 \quad \text{on } \partial\Omega. \end{aligned}$$

Besides the layer at the outflow boundary $x = 1$, there are also two layers at the boundaries $y = 0$ and $y = 1$, see Figure 2.1. The layer at the outflow boundary is often called exponential layer and the layers parallel to the flow direction parabolic layers.

Layers are very small structures: the size of exponential layers is $\mathcal{O}(\varepsilon)$ and the size of parabolic layers $\mathcal{O}(\sqrt{\varepsilon})$. Usually, such structures cannot be resolved on given grids, i.e., they cannot be represented on these grids. Clearly, a structure that cannot be represented cannot be simulated. However, the Galerkin finite element method tries to simulate all important features.

A result, computed with the Galerkin finite element method, for $\varepsilon = 10^{-8}$ and the P_1 finite element method on a grid consisting of $32 \times 32 \times 2$ triangles, which corresponds to 1089 degrees of freedom (including Dirichlet nodes) is presented in Figure 2.2. One obtains a solution that is globally polluted with huge spurious oscillations. The numerical approximation is completely useless. \square

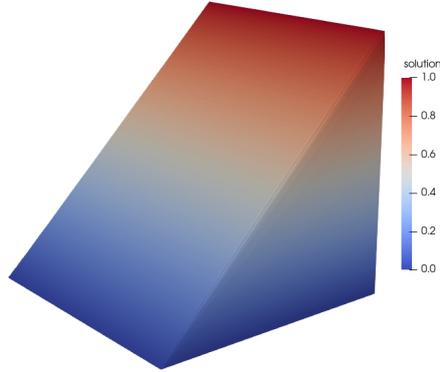


Fig. 2.1 Example 2.5, solution.

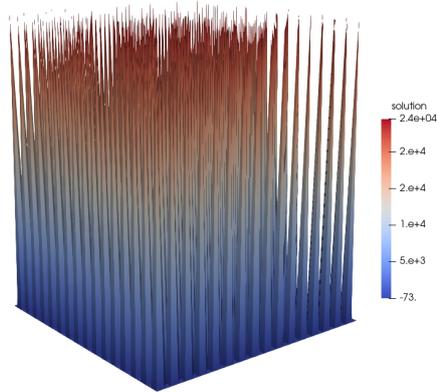


Fig. 2.2 Example 2.5, numerical solution obtained with the Galerkin finite element method, note the size of the values.

2.3 Stabilized Finite Element Methods

Remark 2.6. On the $H^1(\Omega)$ norm for the numerical analysis of convection-dominated problems. Consider the problem: Find $u \in V = H_0^1(\Omega)$ such that

$$a(u, v) = f(v) \quad \forall v \in V \quad (2.3)$$

with

$$a(u, v) := \varepsilon (\nabla u, \nabla v) + (\mathbf{b} \cdot \nabla u, v) + (\sigma u, v), \quad f(v) := (f, v).$$

Let the condition

$$\left(-\frac{1}{2}\nabla \cdot \mathbf{b} + \sigma\right)(\mathbf{x}) \geq \mu_0 > 0 \quad \text{almost everywhere in } \Omega \quad (2.4)$$

be satisfied, which is stronger than condition (1.30). Then, an analogous calculation as in the proof of Theorem 1.11 shows that $a(\cdot, \cdot)$ is uniformly coercive with respect to the following norm, which depends on ε ,

$$\|v\|_\varepsilon^2 := \varepsilon \|\nabla v\|_{L^2(\Omega)}^2 + \mu_0 \|v\|_{L^2(\Omega)}^2,$$

i.e., there is a constant m which does not depend on ε such that

$$a(v, v) \geq m \|v\|_\varepsilon^2 \quad \forall v \in V. \quad (2.5)$$

Applying integration by parts, *exercise*, shows that there is a constant M , which is also independent of ε , such that

$$|a(v, w)| \leq M \|v\|_\varepsilon \|w\|_{H^1(\Omega)} \quad \forall (v, w) \in V \times V. \quad (2.6)$$

However, there is no constant \tilde{M} , which is independent of ε , with

$$|a(v, w)| \leq \tilde{M} \|v\|_\varepsilon \|w\|_\varepsilon \quad \forall (v, w) \in V \times V.$$

Using the estimates (2.5) and (2.6) with constants that are independent of ε , one obtains in a similar way as in the proof of the Lemma of Cea that

$$\|u - u^h\|_\varepsilon \leq C \inf_{v^h \in V^h} \|u - v^h\|_{H^1(\Omega)},$$

with C independent of ε . If V^h is a standard finite element space (piecewise polynomial), then one can show that it is for the best approximation error in layers

$$\inf_{v^h \in V^h} \|u - v^h\|_{H^1(\Omega)} \rightarrow \infty \quad \text{for } \varepsilon \rightarrow 0,$$

for fixed h . Consequently, there is no uniform convergence $\|u - u^h\|_\varepsilon \rightarrow 0$ for $h \rightarrow 0$. The norm $\|\cdot\|_\varepsilon$ is not suited for the investigation of numerical methods for convection-dominated problems. It turns out that the use of appropriate norms is important for the numerical analysis of discretizations for convection-dominated problems. \square

2.3.1 The Streamline-Upwind Petrov–Galerkin (SUPG) Method

Remark 2.7. Goal. The goal consists in the construction of a method that is more stable than the Galerkin finite element method and which can be used with finite elements of arbitrary order. The convergence of this method, in

an appropriate norm, should be of higher order. In addition, the constant in the error estimate should not blow up if $\varepsilon \rightarrow 0$. Such estimates are called robust error estimates.

Consider problem (1.7) and assume that condition (2.4) is satisfied. \square

Remark 2.8. The basic idea. The basic idea consists in a penalization of large values of the so-called strong residual. Such methods are called residual-based stabilizations.

Given a linear partial differential equation in strong form

$$A_{\text{str}} u_{\text{str}} = f, \quad f \in L^2(\Omega),$$

and its Galerkin finite element discretization: Find $u^h \in V^h$ such that

$$a^h(u^h, v^h) = (f, v^h) \quad \forall v^h \in V^h. \quad (2.7)$$

For residual-based stabilizations, a modification of A_{str} is needed which is well-defined for finite element functions. This modification should be also a linear operator and it is denoted by $A_{\text{str}}^h : V^h \rightarrow L^2(\Omega)$. The (strong) residual is now defined by

$$r^h(u^h) = A_{\text{str}}^h u^h - f \in L^2(\Omega).$$

In general, it holds $r^h(u^h) \neq 0$, but a good numerical approximation of the solution of the continuous problem should have in some sense a small residual. Now, instead of finding the solution of (2.7), the minimizer of the residual is searched, i.e, the following optimization problem is considered

$$\arg \min_{u^h \in V^h} \|r^h(u^h)\|_{L^2(\Omega)}^2 = \arg \min_{u^h \in V^h} (r^h(u^h), r^h(u^h)). \quad (2.8)$$

The necessary condition for taking the minimum is the vanishing of the Gâteaux derivative. This derivative is computed by using the linearity of A_{str}^h and the bilinearity of the inner product in $L^2(\Omega)$

$$\begin{aligned} 0 &= \lim_{\varepsilon \rightarrow 0} \frac{(r^h(u^h + \varepsilon v^h), r^h(u^h + \varepsilon v^h)) - (r^h(u^h), r^h(u^h))}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0} \frac{(r^h(u^h) + \varepsilon A_{\text{str}}^h v^h, r^h(u^h) + \varepsilon A_{\text{str}}^h v^h) - (r^h(u^h), r^h(u^h))}{\varepsilon} \\ &= 2(r^h(u^h), A_{\text{str}}^h v^h) \quad \forall v^h \in V^h. \end{aligned}$$

It follows that the necessary condition for the solution of (2.8) is

$$(r^h(u^h), A_{\text{str}}^h v^h) = 0 \quad \forall v^h \in V^h.$$

A generalization consists in considering the minimization problem

$$\arg \min_{u^h \in V^h} \left\| \delta^{1/2} r^h(u^h) \right\|_{L^2(\Omega)}^2 = \arg \min_{u^h \in V^h} (\delta r^h(u^h), r^h(u^h)) \quad (2.9)$$

with the positive weighting function $\delta(\mathbf{x})$. Analogously to the derivation for the special case, one obtains as necessary condition for the minimum

$$(\delta r^h(u^h), A_{\text{str}}^h v^h) = 0 \quad \forall v^h \in V^h. \quad (2.10)$$

The solutions of (2.8) or (2.9) will not be identical to the solution of the Galerkin discretization (2.7). It turns out that the reason for the Galerkin discretization to fail is that the solution possesses structures (scales) that are important but which are not resolved by the used finite element space (grid). For convection-diffusion problems, such structures are layers, e.g., at boundaries. The numerical methods should also compute sharp layers. However the sharpness of layers in numerical solutions is restricted by the resolution of the finite element space, which is generally much coarser than the layer width. Hence, even for a numerical solution with sharp layers, the residual in the layer regions are very large. In particular, a numerical solution with sharp layers (with respect to the resolution of the finite element space) will not be the minimizer of (2.8) or (2.9), see Figure 2.3. The minimizers of (2.8) or (2.9) tend to possess strongly smeared layers and these solutions are useless in applications. For this reason, one considers in residual-based stabilizations a combination of the Galerkin discretization (2.7), which possesses not sufficient diffusion, and the minimization of the residual, which is over-diffusive,

$$a^h(u^h, v^h) + (\delta r^h(u^h), A_{\text{str}}^h v^h) = (f, v^h) \quad \forall v^h \in V^h. \quad (2.11)$$

The goal of numerical analysis consists in determining the weighting function δ optimally in an asymptotic sense. \square

Definition 2.9. Streamline-Upwind Petrov–Galerkin FEM, SUPG method, Streamline-Diffusion FEM, SDFEM. The Streamline-Upwind Petrov–Galerkin (SUPG) FEM or Streamline-Diffusion FEM (SDFEM) has the form: Find $u^h \in V^h$, such that

$$a^h(u^h, v^h) = f^h(v^h) \quad \forall v^h \in V^h \quad (2.12)$$

with $V^h \subset V$ and

$$\begin{aligned} a^h(v, w) &:= a(v, w) & (2.13) \\ &+ \sum_{K \in \mathcal{T}^h} \int_K \delta_K \left(-\varepsilon \Delta v(\mathbf{x}) + \mathbf{b}(\mathbf{x}) \cdot \nabla v(\mathbf{x}) + \sigma(\mathbf{x}) v(\mathbf{x}) \right) \left(\mathbf{b}(\mathbf{x}) \cdot \nabla w(\mathbf{x}) \right) d\mathbf{x}, \\ f^h(w) &:= (f, w) + \sum_{K \in \mathcal{T}^h} \int_K \delta_K f(\mathbf{x}) \left(\mathbf{b}(\mathbf{x}) \cdot \nabla w(\mathbf{x}) \right) d\mathbf{x}. \end{aligned}$$

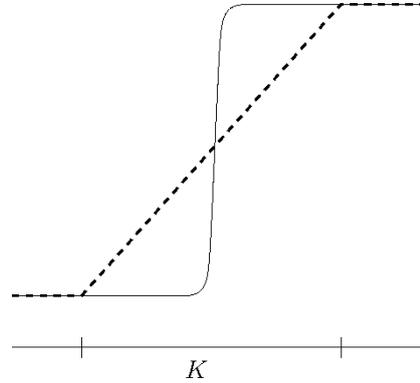


Fig. 2.3 Function with sharp layer (solid line) and optimal piecewise linear approximation in a mesh cell K (dashed line). The equation that is fulfilled by the function in K is far from being satisfied by the piecewise linear approximation. Hence, despite the approximation is of the type considered to be optimal, the residual will be large.

Here, $\{\delta_K\}$ are user-chosen weights, which are called stabilization parameters or SUPG parameters. \square

Remark 2.10. Concerning the SUPG method.

- The method was developed in Hughes & Brooks (1979); Brooks & Hughes (1982).
- The name ‘SUPG’ comes from the fact that the method can be considered as a Petrov–Galerkin method, i.e., a finite element method with different test and ansatz spaces, with the test space

$$\text{span} \left\{ w(\mathbf{x}) + \sum_{K \in \mathcal{T}^h} \delta_K \mathbf{b}(\mathbf{x}) \cdot \nabla w(\mathbf{x}) \right\}.$$

- The SUPG method introduces artificial diffusion only in the so-called streamline direction $\mathbf{b}(\mathbf{x}) \cdot \nabla w(\mathbf{x})$. From this property, the name ‘Streamline Diffusion FEM’ originates.
- The operator A_{str}^h is given in the second part of the bilinear form (2.13). The second derivative for finite element functions is defined only piecewise.
- In the stabilization term of the SUPG method, not the strong operator A_{str}^h applied to the test function is used, as in (2.11), but only the first order term contained in this expression. However, for convection-dominated problems, the first order term is the dominating term of the strong operator applied to the test function. The numerical analysis presented below will show that using only the first order term suffices.
- Altogether, the SUPG method is the most popular stabilized finite element method in academics. However, since there are usually still notable

spurious oscillations in a vicinity of layers, compare Example 2.23 below, it is only of restricted usefulness in practice, e.g., see John *et al.* (2009).

- Generally, the SUPG parameter is a general function. However, in practice it is often chosen as a piecewise constant function. The goal of the finite element error analysis consists in proposing a good asymptotic choice of this parameter.

□

Definition 2.11. Consistent finite element method. Let $u(\mathbf{x})$ be a sufficiently smooth solution of: Find $u \in V$ such that

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

where $a(\cdot, \cdot)$ is an appropriate bilinear form and $f(\cdot)$ an appropriate functional. A finite element method related to this problem: Find $u^h \in V^h$ such that

$$a^h(u^h, v^h) = f^h(v^h) \quad \forall v^h \in V^h$$

is called consistent, if

$$a^h(u, v^h) = f^h(v^h) \quad \forall v^h \in V^h. \quad (2.14)$$

□

Remark 2.12. Consistency. Note that consistency of a finite element method is not the same as consistency of a finite difference method, see the lecture notes on Numerical Mathematics 3. For finite element methods, consistency means that a sufficiently smooth solution satisfies also the discrete equation.

□

Lemma 2.13. Galerkin orthogonality. *A consistent finite element method has the property of the Galerkin orthogonality*

$$a^h(u - u^h, v^h) = 0 \quad \forall v^h \in V^h. \quad (2.15)$$

The error is ‘orthogonal’ to the finite element space.

Proof. The statement of the lemma follows immediately by subtracting (2.12) and (2.14). ■

Lemma 2.14. Consistency of the SUPG method. *The SUPG method (2.12), (2.13) is consistent.*

Proof. A sufficiently smooth solution u of (1.27) satisfies the strong form of the equation even pointwise. Hence, the residual is pointwise zero. Inserting this solution in the SUPG formulation (2.12), (2.13) results in a vanishing of the stabilization term. It remains

$$a(u, v^h) = f(v^h) \quad \forall v^h \in V^h,$$

which is satisfied by the weak solution since $V^h \subset V$. That means, the smooth solution satisfies also the discrete equation. ■

Definition 2.15. SUPG norm. Let for almost all $\mathbf{x} \in \Omega$ condition (2.4) be satisfied. In V^h , the SUPG norm is defined by

$$\begin{aligned} & \|v^h\|_{\text{SUPG}} \\ & := \left(\varepsilon |v^h|_{H^1(\Omega)}^2 + \mu_0 \|v^h\|_{L^2(\Omega)}^2 + \sum_{K \in \mathcal{T}^h} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla v^h) \right\|_{L^2(K)}^2 \right)^{1/2} \end{aligned} \quad (2.16)$$

□

Theorem 2.16 (Coercivity of the SUPG bilinear form). Assume that $\mathbf{b} \in W^{1,\infty}(\Omega)$, $\sigma \in L^\infty(\Omega)$, (2.4), and let for each $K \in \mathcal{T}^h$

$$0 < \delta_K \leq \frac{1}{2} \min \left\{ \frac{h_K^2}{\varepsilon C_{\text{inv}}^2}, \frac{\mu_0}{\|\sigma\|_{L^\infty(K)}} \right\}, \quad (2.17)$$

where C_{inv} is the constant in the inverse estimate (2.2). Then, the SUPG bilinear form is coercive with respect to the SUPG norm, concretely, it is

$$a^h(v^h, v^h) \geq \frac{1}{2} \|v^h\|_{\text{SUPG}}^2 \quad \forall v^h \in V^h.$$

Proof. Integration by parts gives, see the proof of Theorem 1.11,

$$(\mathbf{b} \cdot \nabla v^h + \sigma v^h, v^h) = \left(\left(-\frac{\nabla \cdot \mathbf{b}}{2} + \sigma \right) v^h, v^h \right) \quad \forall v^h \in V^h.$$

With the definition of μ_0 , one obtains

$$\begin{aligned} & a^h(v^h, v^h) \\ & = \varepsilon |v^h|_1^2 + \underbrace{\int_{\Omega} \left(\sigma(\mathbf{x}) - \frac{\nabla \cdot \mathbf{b}(\mathbf{x})}{2} \right) (v^h)^2(\mathbf{x}) \, d\mathbf{x}}_{\geq \mu_0 > 0} + \sum_{K \in \mathcal{T}^h} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla v^h) \right\|_{L^2(K)}^2 \\ & \quad + \sum_{K \in \mathcal{T}^h} \int_K \delta_K (-\varepsilon \Delta v^h(\mathbf{x}) + \sigma(\mathbf{x}) v^h(\mathbf{x})) (\mathbf{b}(\mathbf{x}) \cdot \nabla v^h(\mathbf{x})) \, d\mathbf{x} \quad (2.18) \\ & \geq \|v^h\|_{\text{SUPG}}^2 - \left| \sum_{K \in \mathcal{T}^h} \int_K \delta_K (-\varepsilon \Delta v^h(\mathbf{x}) + \sigma(\mathbf{x}) v^h(\mathbf{x})) (\mathbf{b}(\mathbf{x}) \cdot \nabla v^h(\mathbf{x})) \, d\mathbf{x} \right|. \end{aligned}$$

Now, the last term will be estimated from above. Then, one obtains altogether an estimate from below if the estimate of the last term is subtracted from

the first term. In the following estimate, one uses the conditions (2.17) on the SUPG parameter. It is for each $K \in \mathcal{T}^h$

$$\begin{aligned}
& \left| \int_K \delta_K (-\varepsilon \Delta v^h(\mathbf{x}) + \sigma(\mathbf{x}) v^h(\mathbf{x})) (\mathbf{b} \cdot \nabla v^h(\mathbf{x})) \, d\mathbf{x} \right| \\
& \leq \int_K \left(\delta_K^{1/2} \varepsilon |\Delta v^h(\mathbf{x})| \right) \left(\delta_K^{1/2} |\mathbf{b} \cdot \nabla v^h(\mathbf{x})| \right) \, d\mathbf{x} \\
& \quad + \int_K \left(\delta_K^{1/2} |\sigma(\mathbf{x})| |v^h(\mathbf{x})| \right) \left(\delta_K^{1/2} |\mathbf{b} \cdot \nabla v^h(\mathbf{x})| \right) \, d\mathbf{x} \\
& \stackrel{\text{CS}}{\leq} \left(\delta_K^{1/2} \varepsilon \|\Delta v^h\|_{L^2(K)} + \delta_K^{1/2} \|\sigma\|_{L^\infty(K)} \|v^h\|_{L^2(K)} \right) \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla v^h) \right\|_{L^2(K)} \\
& \stackrel{(2.2)}{\leq} \left(\delta_K^{1/2} \frac{\varepsilon C_{\text{inv}}}{h_K} \|\nabla v^h\|_{L^2(K)} + \delta_K^{1/2} \|\sigma\|_{L^\infty(K)} \|v^h\|_{L^2(K)} \right) \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla v^h) \right\|_{L^2(K)} \\
& \stackrel{(2.17)}{\leq} \left(\frac{h_K}{\sqrt{2\varepsilon} C_{\text{inv}}} \frac{\varepsilon C_{\text{inv}}}{h_K} \|\nabla v^h\|_{L^2(K)} + \frac{\sqrt{\mu_0}}{\sqrt{2} \|\sigma\|_{L^\infty(K)}} \|\sigma\|_{L^\infty(K)} \|v^h\|_{L^2(K)} \right) \\
& \quad \times \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla v^h) \right\|_{L^2(K)} \\
& = \left(\sqrt{\frac{\varepsilon}{2}} \|\nabla v^h\|_{L^2(K)} + \sqrt{\frac{\mu_0}{2}} \|v^h\|_{L^2(K)} \right) \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla v^h) \right\|_{L^2(K)} \\
& \stackrel{\text{Young}}{\leq} \frac{\varepsilon}{2} \|\nabla v^h\|_{L^2(K)}^2 + \frac{1}{4} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla v^h) \right\|_{L^2(K)}^2 + \frac{\mu_0}{2} \|v^h\|_{L^2(K)}^2 \\
& \quad + \frac{1}{4} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla v^h) \right\|_{L^2(K)}^2 \\
& = \frac{1}{2} \|v^h\|_{\text{SUPG}(K)}^2.
\end{aligned}$$

Now, the proof is finished by summing over all mesh cells and inserting the result in (2.18). \blacksquare

Corollary 2.17 (Coercivity of the SUPG bilinear form for P_1 finite elements). *Let the assumptions of Theorem 2.16 with respect to the coefficients of the problem be satisfied. For piecewise linear finite elements, the SUPG bilinear form (2.13) is coercive with respect to the SUPG norm if*

$$0 < \delta_K \leq \frac{\mu_0}{\|\sigma\|_{L^\infty(K)}^2}. \quad (2.19)$$

Proof. The proof is the same as for Theorem 2.16, where one uses that for piecewise linear finite elements $\Delta v^h(\mathbf{x})|_K = 0$ for all $K \in \mathcal{T}^h$. Thus, the corresponding terms do not appear in the proof. \blacksquare

Corollary 2.18 (Existenz and uniqueness of a solution of the SUPG method). *Let the assumptions of Theorem 2.16 and Corollary 2.17, respec-*

tively, be satisfied. Then, the SUPG finite element method (2.12), (2.13) has a unique solution.

Proof. The statement is obtained by the application of the Theorem of Lax–Milgram. The coercivity of the bilinear form was proved in Theorem 2.16 and Corollary 2.17, respectively. For the boundedness, one uses similar estimates as in the proof of Theorems 2.16 and 1.11, *exercise*. ■

Remark 2.19. On the coercivity of the SUPG bilinear form.

- The proof of Theorem 2.16 is typical for the numerical analysis of stabilized finite element methods. One tries to get rid of the troubling terms by estimating them with the used norm. This approach works only if one uses an appropriate norm. In particular, the stabilization has to appear in the norm.
- Theorem 2.16 provides an upper bound for the SUPG parameter. This bound is generally not critical in applications.
- From Theorem 2.16, one obtains the stability of the SUPG method with respect to the SUPG norm. Stability means that an appropriate norm of the solution can be estimated with the data of the problem. It is

$$\begin{aligned}
& \|u^h\|_{\text{SUPG}}^2 \\
& \leq 2a^h(u^h, u^h) = 2f^h(u^h) \\
& = 2(f, u^h) + 2 \sum_{K \in \mathcal{T}^h} \int_K \delta_K f(\mathbf{x}) (\mathbf{b}(\mathbf{x}) \cdot \nabla u^h(\mathbf{x})) \, d\mathbf{x} \\
& \stackrel{\text{CS}}{\leq} \frac{2}{\sqrt{\mu_0}} \|f\|_{L^2(\Omega)} \sqrt{\mu_0} \|u^h\|_{L^2(\Omega)} \\
& \quad + 2 \sum_{K \in \mathcal{T}^h} \left\| \delta_K^{1/2} f \right\|_{L^2(K)} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla u^h) \right\|_{L^2(K)} \\
& \stackrel{\text{Young}}{\leq} C \|f\|_{L^2(\Omega)}^2 + \frac{1}{2} \left(\mu_0 \|u^h\|_{L^2(\Omega)}^2 + \sum_{K \in \mathcal{T}^h} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla u^h) \right\|_{L^2(K)}^2 \right) \\
& \leq C \|f\|_{L^2(\Omega)}^2 + \frac{1}{2} \|u^h\|_{\text{SUPG}}^2.
\end{aligned}$$

In the last estimate, it was used that the terms are part of the SUPG norm and thus they can be estimated by the whole norm. Now, the second term on the right-hand side can be absorbed in the left-hand side and stability is proved. The stability constant depends on μ_0 and on the upper bounds of δ_K .

- All $v^h \in V^h$ satisfy

$$\|v^h\|_{\text{SUPG}} \geq \min \{1, \mu_0^{1/2}\} \|v^h\|_\varepsilon.$$

Hence, the SUPG method is also stable with respect to the norm $\|\cdot\|_\varepsilon$. With respect to this norm, also the Galerkin finite element method is stable, however this method is not stable with respect to the SUPG norm. That means, the stability of the SUPG method is stronger than the stability of the Galerkin finite element method. \square

Theorem 2.20 (Convergence of the SUPG method). *Let the solution of (1.27) satisfy $u \in H^{k+1}(\Omega)$, $k \geq 1$, let $\mathbf{b} \in W^{1,\infty}(\Omega)$, $\sigma \in L^\infty(\Omega)$, let the assumptions of Theorem 2.16 be satisfied, and consider the SUPG method for P_k finite elements, $k \geq 1$. Let the SUPG parameter be given as follows*

$$\delta_K = \begin{cases} C_0 \frac{h_K^2}{\varepsilon} & \text{for } h_K < \varepsilon, \\ C_0 h_K & \text{for } \varepsilon \leq h_K, \end{cases} \quad (2.20)$$

where the constant $C_0 > 0$ is sufficiently small such that (2.17) is satisfied for $k \geq 2$ or (2.19) for $k = 1$, respectively. Then, the solution $u^h \in P_k$ of the SUPG method (2.12), (2.13) satisfies the following error estimate

$$\|u - u^h\|_{\text{SUPG}} \leq C \left(\varepsilon^{1/2} h^k + h^{k+1/2} \right) |u|_{H^{k+1}(\Omega)},$$

where the constant C does not depend on ε and h .

Proof. Let $u_I^h \in V^h$ be the Lagrangian interpolant of $u(\mathbf{x})$. One obtains with the triangle inequality

$$\|u - u^h\|_{\text{SUPG}} \leq \|u - u_I^h\|_{\text{SUPG}} + \|u_I^h - u^h\|_{\text{SUPG}}. \quad (2.21)$$

The first term on the right-hand side is the interpolation error. Note that for both regimes, it is

$$\delta_K \leq C_0 h_K \leq Ch. \quad (2.22)$$

Using this property after having applied the interpolation error estimate (2.1) to each term of the SUPG norm individually gives

$$\begin{aligned} \|u - u_I^h\|_{\text{SUPG}} &\leq \left(C\varepsilon h^{2k} |u|_{H^{k+1}(\Omega)}^2 + C\mu_0 h^{2(k+1)} |u|_{H^{k+1}(\Omega)}^2 \right. \\ &\quad \left. + C \sum_{K \in \mathcal{T}^h} \delta_K \|\mathbf{b}\|_{L^\infty(K)}^2 h_K^{2k} |u|_{H^{k+1}(K)}^2 \right)^{1/2} \\ &\leq C \left(\varepsilon h^{2k} + h^{2(k+1)} + h^{2k+1} \right)^{1/2} |u|_{H^{k+1}(\Omega)} \\ &\leq C \left(\varepsilon^{1/2} h^k + h^{k+1/2} \right) |u|_{H^{k+1}(\Omega)}. \end{aligned}$$

Consider now the second term on the right-hand side of (2.21). The coercivity, Theorem 2.16, and the Galerkin orthogonality yield

$$\frac{1}{2} \|u_I^h - u^h\|_{\text{SUPG}}^2 \leq a^h(u_I^h - u^h, u_I^h - u^h) = a^h(u_I^h - u, u_I^h - u^h).$$

Now, the triangle inequality is applied to $a^h(u_I^h - u, u_I^h - u^h)$ and every term is bounded individually. In these estimates, the interpolation estimate (2.1) plays an important role. Let $w^h = u_I^h - u^h$. One obtains for the diffusive term

$$\begin{aligned} & |\varepsilon (\nabla(u_I^h - u), \nabla w^h)| \\ & \stackrel{\text{CS}}{\leq} \varepsilon \|\nabla(u_I^h - u)\|_{L^2(\Omega)} \|\nabla w^h\|_{L^2(\Omega)} = \varepsilon^{1/2} \|\nabla(u_I^h - u)\|_{L^2(\Omega)} \varepsilon^{1/2} \|\nabla w^h\|_{L^2(\Omega)} \\ & \stackrel{(2.1)}{\leq} C \varepsilon^{1/2} h^k |u|_{H^{k+1}(\Omega)} \varepsilon^{1/2} \|\nabla w^h\|_{L^2(\Omega)} \leq C \varepsilon^{1/2} h^k |u|_{H^{k+1}(\Omega)} \|w^h\|_{\text{SUPG}}. \end{aligned}$$

For the reactive term, one calculates in a similar way

$$\begin{aligned} |(\sigma(u_I^h - u), w^h)| & \stackrel{\text{Hölder}}{\leq} \|\sigma\|_{L^\infty(\Omega)} \|u_I^h - u\|_{L^2(\Omega)} \|w^h\|_{L^2(\Omega)} \\ & = \mu_0^{-1/2} \|\sigma\|_{L^\infty(\Omega)} \|u_I^h - u\|_{L^2(\Omega)} \mu_0^{1/2} \|w^h\|_{L^2(\Omega)} \\ & \stackrel{(2.1)}{\leq} C h^{k+1} |u|_{H^{k+1}(\Omega)} \|w^h\|_{\text{SUPG}}. \end{aligned}$$

Next, the terms are considered that come from the SUPG stabilization. Since for both regimes it is

$$\varepsilon \delta_K \leq C_0 h_K^2,$$

one gets

$$\begin{aligned}
& \left| \sum_{K \in \mathcal{T}^h} (-\varepsilon \Delta (u_I^h - u), \delta_K \mathbf{b} \cdot \nabla w^h)_K \right| \\
& \stackrel{\text{CS}}{\leq} \sum_{K \in \mathcal{T}^h} \varepsilon^{1/2} \|\Delta (u_I^h - u)\|_{L^2(K)} \varepsilon^{1/2} \delta_K^{1/2} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla w^h) \right\|_{L^2(K)} \\
& \leq C_0^{1/2} \sum_{K \in \mathcal{T}^h} h_K \varepsilon^{1/2} \|\Delta (u_I^h - u)\|_{L^2(K)} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla w^h) \right\|_{L^2(K)} \\
& \stackrel{\text{CS}}{\leq} C_0^{1/2} \varepsilon^{1/2} h \left(\sum_{K \in \mathcal{T}^h} \|\Delta (u_I^h - u)\|_{L^2(K)}^2 \right)^{1/2} \\
& \quad \times \left(\sum_{K \in \mathcal{T}^h} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla w^h) \right\|_{L^2(K)}^2 \right)^{1/2} \\
& \stackrel{(2.1)}{\leq} C \varepsilon^{1/2} h \left(\sum_{K \in \mathcal{T}^h} h_K^{2(k-1)} |u|_{H^{k+1}(K)}^2 \right)^{1/2} \left(\sum_{K \in \mathcal{T}^h} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla w^h) \right\|_{L^2(K)}^2 \right)^{1/2} \\
& \leq C \varepsilon^{1/2} h^k |u|_{H^{k+1}(\Omega)} \|w^h\|_{\text{SUPG}}.
\end{aligned}$$

For the other terms, one obtains with (2.22), which holds for both regimes,

$$\begin{aligned}
& \left| \sum_{K \in \mathcal{T}^h} (\mathbf{b} \cdot \nabla (u_I^h - u) + \sigma (u_I^h - u), \delta_K (\mathbf{b} \cdot \nabla w^h))_K \right| \\
& \stackrel{\text{Hölder}}{\leq} \sum_{K \in \mathcal{T}^h} \|\mathbf{b}\|_{L^\infty(K)} \|\nabla (u_I^h - u)\|_{L^2(K)} \delta_K^{1/2} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla w^h) \right\|_{L^2(K)} \\
& \quad + \sum_{K \in \mathcal{T}^h} \|\sigma\|_{L^\infty(K)} \|u_I^h - u\|_{L^2(K)} \delta_K^{1/2} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla w^h) \right\|_{L^2(K)} \\
& \leq C \left(\sum_{K \in \mathcal{T}^h} h_K^{1/2} \|\nabla (u_I^h - u)\|_{L^2(K)} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla w^h) \right\|_{L^2(K)} \right. \\
& \quad \left. + \sum_{K \in \mathcal{T}^h} h_K^{1/2} \|u_I^h - u\|_{L^2(K)} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla w^h) \right\|_{L^2(K)} \right) \\
& \stackrel{\text{CS}}{\leq} Ch^{1/2} \left[\left(\sum_{K \in \mathcal{T}^h} \|\nabla (u_I^h - u)\|_{L^2(K)}^2 \right)^{1/2} + \left(\sum_{K \in \mathcal{T}^h} \|u_I^h - u\|_{L^2(K)}^2 \right)^{1/2} \right] \\
& \quad \times \left(\sum_{K \in \mathcal{T}^h} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla w^h) \right\|_{L^2(K)}^2 \right)^{1/2} \\
& \stackrel{(2.1)}{\leq} C (h^{k+1/2} + h^{k+3/2}) |u|_{H^{k+1}(\Omega)} \|w^h\|_{\text{SUPG}}.
\end{aligned}$$

To obtain an optimal estimate for the convective term, one has to apply first integration by parts

$$\begin{aligned}
(\mathbf{b} \cdot \nabla (u_I^h - u), w^h) &= (\nabla (u_I^h - u), \mathbf{b} w^h) = - (u_I^h - u, \nabla \cdot (\mathbf{b} w^h)) \\
&= - (u_I^h - u, (\nabla \cdot \mathbf{b}) w^h) - (u_I^h - u, \mathbf{b} \cdot \nabla w^h).
\end{aligned}$$

Both terms on the right-hand side are bounded separately. Using the same tools as for the other estimates yields

$$\begin{aligned}
& |(u_I^h - u, (\nabla \cdot \mathbf{b}) w^h)| \\
& \leq \mu_0^{-1/2} \|\nabla \cdot \mathbf{b}\|_{L^\infty(\Omega)} \left(\sum_{K \in \mathcal{T}^h} \|u_I^h - u\|_{L^2(K)}^2 \right)^{1/2} \mu_0^{1/2} \|w^h\|_{L^2(\Omega)} \\
& \leq Ch^{k+1} |u|_{H^{k+1}(\Omega)} \|w^h\|_{\text{SUPG}}.
\end{aligned}$$

In the estimate of the other term, one has to distinguish whether in the mesh cell K it is $\varepsilon \leq h_K$ or $\varepsilon > h_K$. One gets

$$\begin{aligned}
& |(u_I^h - u, \mathbf{b} \cdot \nabla w^h)| \\
& \stackrel{\text{H\"older}}{\leq} \sum_{\varepsilon \leq h_K} \delta_K^{-1/2} \|u_I^h - u\|_{L^2(K)} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla w^h) \right\|_{L^2(K)} \\
& \quad + \sum_{\varepsilon > h_K} \|\mathbf{b}\|_{L^\infty(K)} \|u_I^h - u\|_{L^2(K)} \|\nabla w^h\|_{L^2(K)} \\
& \stackrel{(2.1)}{\leq} C \left(\sum_{\varepsilon \leq h_K} \delta_K^{-1/2} h_K^{k+1} |u|_{H^{k+1}(K)} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla w^h) \right\|_{L^2(K)} \right. \\
& \quad \left. + \sum_{\varepsilon > h_K} h_K^{k+1} |u|_{H^{k+1}(K)} \|\nabla w^h\|_{L^2(K)} \right) \\
& \stackrel{C_0 h_K = \delta_K, \varepsilon > h_K}{\leq} C \left(\sum_{\varepsilon \leq h_K} C_0^{-1/2} h_K^{-1/2} h_K^{k+1} |u|_{H^{k+1}(K)} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla w^h) \right\|_{L^2(K)} \right. \\
& \quad \left. + \sum_{\varepsilon > h_K} h_K^{k+1/2} |u|_{H^{k+1}(K)} \varepsilon^{1/2} \|\nabla w^h\|_{L^2(K)} \right) \\
& \stackrel{\text{CS}}{\leq} C h^{k+1/2} |u|_{H^{k+1}(\Omega)} \left[\left(\sum_{K \in \mathcal{T}^h} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla w^h) \right\|_{L^2(K)}^2 \right)^{1/2} + \varepsilon^{1/2} |w^h|_1 \right] \\
& \leq C h^{k+1/2} |u|_{H^{k+1}(\Omega)} \|w^h\|_{\text{SUPG}}.
\end{aligned}$$

Summarizing all estimates, the statement of the theorem is proved. \blacksquare

Remark 2.21. Concerning the error estimate.

- In the convection-dominated regime, where $\varepsilon \ll h$, the order of error reduction in the SUPG norm is $k + 1/2$ and in the diffusion-dominated case, the order of convergence is k . In the latter case, the SUPG norm is essentially the $H^1(\Omega)$ semi norm such that order k is optimal.
- It is essential for obtaining an estimate with a constant C which is independent on inverse powers of ε that the term

$$\left(\sum_{K \in \mathcal{T}^h} \left\| \delta_K^{1/2} (\mathbf{b} \cdot \nabla w^h) \right\|_{L^2(K)}^2 \right)^{1/2}$$

is part of the norm, which is used for estimating the error. Such a robust estimate does not hold for the norm $\|\cdot\|_\varepsilon$.

- For the interpretation of the results, one has to take into account that different stabilization parameters by choosing different values of C_0 lead also to different norms on the left-hand side of the estimate.

- On the other hand, the practical importance of a constant which is independent of ε is somewhat questionable since in general $|u|_{H^{k+1}(\Omega)}$ depends on ε .
- In numerical simulations for convection-dominated problems, often one can observe even a reduction of order h^{k+1} for the error in $L^2(\Omega)$, in particular on structured grids. However, in Zhou (1997) examples were constructed that show that the estimate of Theorem 2.20 is sharp also for the error in $L^2(\Omega)$.

□

Remark 2.22. Different choices of the SUPG parameter.

- A refined analysis, taking the polynomial degree k of the finite element into account, proposes the stabilization parameter

$$\delta_K = \begin{cases} C_0 \frac{h_K}{\|\mathbf{b}\|_{L^\infty(K)}} & \text{for } \text{Pe}_K \geq 1, \\ C_0 \frac{h_K^2}{\varepsilon} & \text{else,} \end{cases} \quad \text{with } \text{Pe}_K = \frac{\|\mathbf{b}\|_{L^\infty(K)} h_K}{2k\varepsilon}. \quad (2.23)$$

- In practice, one takes for linear and d -linear finite elements instead of (2.20) also the parameter

$$\delta_K = \frac{h_K}{2\|\mathbf{b}\|_{L^\infty(K)}} \left(\coth(\text{Pe}_K) - \frac{1}{\text{Pe}_K} \right), \quad \text{Pe}_K = \frac{\|\mathbf{b}\|_{L^\infty(K)} h_K}{2\varepsilon}, \quad (2.24)$$

where Pe_K is the local Péclet number, since in one dimensions one recovers under certain conditions the Iljin–Allen–Southwell scheme, see Definition 3.18. There is no user-chosen constant in this parameter. Asymptotically, both parameters (2.20) and (2.24) have the same behavior.

□

Example 2.23 (SUPG method). The same problem as in Example 2.5 is considered. The solution computed with the SUPG method on the same coarse grid as with the Galerkin method is presented in Figure 2.4. One can see very well that it is much better than the solution obtained with the Galerkin method. However, there are still spurious oscillations, in particular at the parabolic layers at $y = 0$ and $y = 1$. These oscillations are a typical feature of solutions obtained with the SUPG method. They might become smaller with higher order elements or on finer grids. But they will generally vanish only if the layer is resolved.

□

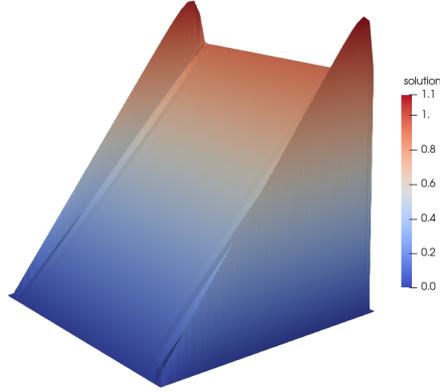


Fig. 2.4 Example 2.23 Result obtained with the SUPG method and the SUPG parameter (2.24).

2.3.2 Other Stabilized Finite Element Methods

Remark 2.24 (Galerkin Least Squares (GLS) method). The GLS method is also a residual-based stabilization. In contrast to the SUPG method, it uses the full linear operator as test function and not only the first order term. Let

$$Lu = -\varepsilon \Delta u + \mathbf{b} \cdot \nabla u + \sigma u,$$

then the GLS method reads as follows: Find $u^h \in V^h$, such that

$$a^h(u^h, v^h) = f^h(v^h) \quad \forall v^h \in V^h$$

with $V^h \subset V$ and

$$\begin{aligned} a^h(v, w) &:= a(v, w) + \sum_{K \in \mathcal{T}^h} \delta_K (Lv, Lw)_K, \\ f^h(w) &:= (f, w) + \sum_{K \in \mathcal{T}^h} \delta_K (f, Lw)_K, \end{aligned}$$

with $\{\delta_K\}$ being the stabilization parameters. In contrast to the SUPG method, the stabilization term of the GLS method is symmetric.

In practice, the GLS method behaves often very similar to the SUPG method. Since the SUPG method requires less terms to be computed, it is usually preferred. \square

Remark 2.25 (Local Projection Stabilization (LPS) method). Another stabilized method with symmetric stabilization term is the LPS method given by: Find $u^h \in V^h$, such that

$$a^h(u^h, v^h) = (f, v^h) \quad \forall v^h \in V^h$$

with $V^h \subset V$ and

$$a^h(v, w) := a(v, w) + \sum_{M \in \mathcal{M}^h} \delta_M (\kappa^h(\mathbf{b} \cdot \nabla v), \kappa^h(\mathbf{b} \cdot \nabla w))_M. \quad (2.25)$$

Here M are so-called macro cells, e.g., patches of mesh cells. Then, a finite element space $D^h(M)$ is defined on the macro cells and a local projection operator $\pi^h : L^2(\Omega) \rightarrow D^h = \cup D^h(M)$. The so-called fluctuation operator is defined by $\kappa^h = \text{id} - \pi^h$. Furthermore, δ_M are stabilization parameters that have to be chosen appropriately.

The stabilization term (2.25) is, apart of the fluctuation operator, the main term of the SUPG stabilization. The local projection operator maps to a large scale finite element space. Then, the fluctuation operator is the difference of all scales of the finite element functions and the large scales, hence it can be interpreted to represent the small scales of the finite element function. Thus, (2.25) adds only additional diffusion to small scales of the numerical solution. By using the fluctuations and an appropriate choice of the stabilization parameters, the consistency error of this method becomes of sufficiently high order.

In contrast to the SUPG and GLS methods, there is no modification of the right-hand side in LPS methods. However, the LPS method leads to a larger matrix stencil. In numerical studies it was shown that also the numerical approximations with the LPS method usually exhibit spurious oscillations in a vicinity of layers. \square

Remark 2.26 (Continuous Interior Penalty (CIP) method, edge stabilization). It was found that stability with respect to dominating convection can be obtained also by a different mechanism, namely by penalizing jumps of the gradient of finite element functions. For the Galerkin FEM, there are huge spurious oscillations and consequently huge jumps of the gradient. This method is called CIP method or edge stabilization. It reads as follows: Find $u^h \in V^h$, such that

$$a^h(u^h, v^h) = (f, v^h) \quad \forall v^h \in V^h$$

with $V^h \subset V$ and

$$a^h(v, w) := a(v, w) + \sum_{F \in \mathcal{F}^h} \delta_F (\mathbf{b} \cdot \llbracket \nabla v \rrbracket_F, \mathbf{b} \cdot \llbracket \nabla w \rrbracket_F),$$

where \mathcal{F}^h is the set of all interior facets (edges in 2d), $\llbracket \nabla v \rrbracket_F$ is the jump of ∇v across F , where for each F an arbitrary but fixed normal vector has to be chosen to fix the direction of the jump, $(\cdot, \cdot)_F$ is the inner product in $L^2(F)$, and δ_F is a stabilization parameter.

This method is a symmetric stabilization method and it does not change the right-hand side. It also enlarges the matrix stencil compared with the

Galerkin discretization or the SUPG method. Numerical solutions computed with the CIP method possess usually spurious oscillations in a vicinity of layers. \square

Remark 2.27 (Spurious Oscillations at Layers Diminishing (SOLD) methods). All stabilized methods introduced so far compute numerical solutions with notable spurious oscillations in a vicinity of layers. Such oscillations are often unacceptable in practice. There have been many proposals in the literature to extend these methods with additional terms that should suppress the spurious oscillations. The class of these methods is called SOLD methods. Often, the additional terms are nonlinear, since a corresponding stabilization parameter depends on the numerical solution. This approach seems to be reasonable because the solution behaves completely different at layers and away from layers.

However, it turned out in numerical studies that there is no SOLD method that really removes the spurious oscillations, e.g., compare John & Knobloch (2007). Meanwhile, there are other nonlinear methods, which do not rely on traditional stabilized finite element methods, that achieve the goal of computing oscillation-free and accurate numerical solutions in many situations. \square