

II.1. A residual error estimator for the model problem

II.1.1. The model problem. As a model problem we consider the Poisson equation with mixed Dirichlet-Neumann boundary conditions

$$(II.1.1) \quad \begin{cases} -\Delta u = f & \text{in } \Omega \\ u = 0 & \text{on } \Gamma_D \\ \frac{\partial u}{\partial n} = g & \text{on } \Gamma_N \end{cases}$$

in a connected, bounded, polygonal domain  $\Omega \subset \mathbb{R}^2$  with boundary  $\Gamma$  consisting of two disjoint parts  $\Gamma_D$  and  $\Gamma_N$ . We assume that the Dirichlet boundary  $\Gamma_D$  is closed relative to  $\Gamma$  and has a positive length and that  $f$  and  $g$  are square integrable functions on  $\Omega$  and  $\Gamma_N$ , respectively. The Neumann boundary  $\Gamma_N$  may be empty.

II.1.2. Variational formulation. The standard weak formulation of problem (II.1.1) is:

$$(II.1.2) \quad \begin{cases} \text{Find } u \in H_D^1(\Omega) \text{ such that} \\ \int_{\Omega} \nabla u \cdot \nabla v = \int_{\Omega} f v + \int_{\Gamma_N} g v \\ \text{for all } v \in H_D^1(\Omega). \end{cases}$$

It is well-known that problem (II.1.2) admits a unique solution.

II.1.3. Finite element discretization. We choose an affine equivalent, admissible and shape-regular partition  $\mathcal{T}$  of  $\Omega$  as in Section I.2.7 (p. 14) and consider the following finite element discretization of problem (II.1.2):

$$(II.1.3) \quad \begin{cases} \text{Find } u_{\mathcal{T}} \in S_D^{1,0}(\mathcal{T}) \text{ such that} \\ \int_{\Omega} \nabla u_{\mathcal{T}} \cdot \nabla v_{\mathcal{T}} = \int_{\Omega} f v_{\mathcal{T}} + \int_{\Gamma_N} g v_{\mathcal{T}} \\ \text{for all } v_{\mathcal{T}} \in S_D^{1,0}(\mathcal{T}). \end{cases}$$

Again it is well-known that problem (II.1.3) admits a unique solution.

II.1.4. Equivalence of error and residual. In what follows we always denote by  $u \in H_D^1(\Omega)$  and  $u_{\mathcal{T}} \in S_D^{1,0}(\mathcal{T})$  the exact solutions of problems (II.1.2) and (II.1.3), respectively. They satisfy the identity

$$\int_{\Omega} \nabla(u - u_{\mathcal{T}}) \cdot \nabla v = \int_{\Omega} f v + \int_{\Gamma_N} g v - \int_{\Omega} \nabla u_{\mathcal{T}} \cdot \nabla v$$

for all  $v \in H_D^1(\Omega)$ . The right-hand side of this equation implicitly defines the residual of  $u_{\mathcal{T}}$  as an element of the dual space of  $H_D^1(\Omega)$ .

The Friedrichs and Cauchy-Schwarz inequalities imply for all  $v \in H_D^1(\Omega)$

$$\frac{1}{\sqrt{1 + c_{\Omega}^2}} \|v\|_{H^1(\Omega)} \leq \sup_{\substack{w \in H_D^1(\Omega) \\ \|w\|_{H^1(\Omega)}=1}} \int_{\Omega} \nabla v \cdot \nabla w \leq \|v\|_{H^1(\Omega)}.$$

This corresponds to the fact that the bilinear form

$$H_D^1(\Omega) \ni v, w \mapsto \int_{\Omega} \nabla v \cdot \nabla w$$

defines an isomorphism of  $H_D^1(\Omega)$  onto its dual space. The constants multiplying the first and last term in this inequality are related to the norm of this isomorphism and of its inverse.

The definition of the residual and the above inequality imply the estimate

$$\begin{aligned} & \sup_{\substack{w \in H_D^1(\Omega) \\ \|w\|_{H^1(\Omega)}=1}} \left\{ \int_{\Omega} f w + \int_{\Gamma_N} g w - \int_{\Omega} \nabla u_{\mathcal{T}} \cdot \nabla w \right\} \\ & \leq \|u - u_{\mathcal{T}}\|_{H^1(\Omega)} \\ & \leq \sqrt{1 + c_{\Omega}^2} \sup_{\substack{w \in H_D^1(\Omega) \\ \|w\|_{H^1(\Omega)}=1}} \left\{ \int_{\Omega} f w + \int_{\Gamma_N} g w - \int_{\Omega} \nabla u_{\mathcal{T}} \cdot \nabla w \right\}. \end{aligned}$$

Since the sup-term in this inequality is equivalent to the norm of the residual in the dual space of  $H_D^1(\Omega)$ , we have proved:

$$\text{The norm in } H_D^1(\Omega) \text{ of the error is, up to multiplicative constants, bounded from above and from below by the norm of the residual in the dual space of } H_D^1(\Omega).$$

Most a posteriori error estimators try to estimate this dual norm of the residual by quantities that can more easily be computed from  $f$ ,  $g$ , and  $u_{\mathcal{T}}$ .

II.1.5. Galerkin orthogonality. Since  $S_D^{1,0}(\mathcal{T}) \subset H_D^1(\Omega)$ , the error is orthogonal to  $S_D^{1,0}(\mathcal{T})$ :

$$\int_{\Omega} \nabla(u - u_{\mathcal{T}}) \cdot \nabla w_{\mathcal{T}} = 0$$

for all  $w_{\mathcal{T}} \in S_D^{1,0}(\mathcal{T})$ . Using the definition of the residual, this can be written as

$$\int_{\Omega} f w_{\mathcal{T}} + \int_{\Gamma_N} g w_{\mathcal{T}} - \int_{\Omega} \nabla u_{\mathcal{T}} \cdot \nabla w_{\mathcal{T}} = 0$$

for all  $w_{\mathcal{T}} \in S_D^{1,0}(\mathcal{T})$ . This identity reflects the fact that the discretization (II.1.3) is consistent and that no additional errors are introduced by numerical integration or by inexact solution of the discrete problem. It is often referred to as *Galerkin orthogonality*.

**II.1.6.  $L^2$ -representation of the residual.** Integration by parts element-wise yields for all  $w \in H_D^1(\Omega)$

$$\begin{aligned} & \int_{\Omega} f w + \int_{\Gamma_N} g w - \int_{\Omega} \nabla u_{\mathcal{T}} \cdot \nabla w \\ &= \int_{\Omega} f w + \int_{\Gamma_N} g w - \sum_{K \in \mathcal{T}} \int_K \nabla u_{\mathcal{T}} \cdot \nabla w \\ &= \int_{\Omega} f w + \int_{\Gamma_N} g w + \sum_{K \in \mathcal{T}} \left\{ \int_K \Delta u_{\mathcal{T}} w - \int_{\partial K} \mathbf{n}_K \cdot \nabla u_{\mathcal{T}} w \right\} \\ &= \sum_{K \in \mathcal{T}} \int_K (f + \Delta u_{\mathcal{T}}) w + \sum_{E \in \mathcal{E}_{\Gamma_N}} \int_E (g - \mathbf{n}_E \cdot \nabla u_{\mathcal{T}}) w \\ &\quad - \sum_{E \in \mathcal{E}_{\Omega}} \int_E \mathbb{J}_E(\mathbf{n}_E \cdot \nabla u_{\mathcal{T}}) w. \end{aligned}$$

Here,  $\mathbf{n}_K$  denotes the unit exterior normal to the element  $K$ . Note that  $\Delta u_{\mathcal{T}}$  vanishes on all triangles.

For abbreviation, we define *element* and *edge residuals* by

$$R_K(u_{\mathcal{T}}) = f + \Delta u_{\mathcal{T}}$$

and

$$R_E(u_{\mathcal{T}}) = \begin{cases} -\mathbb{J}_E(\mathbf{n}_E \cdot \nabla u_{\mathcal{T}}) & \text{if } E \in \mathcal{E}_{\Omega}, \\ g - \mathbf{n}_E \cdot \nabla u_{\mathcal{T}} & \text{if } E \in \mathcal{E}_{\Gamma_N}, \\ 0 & \text{if } E \in \mathcal{E}_{\Gamma_D}. \end{cases}$$

Then we obtain the following  $L^2$ -representation of the residual

$$\begin{aligned} & \int_{\Omega} f w + \int_{\Gamma_N} g w - \int_{\Omega} \nabla u_{\mathcal{T}} \cdot \nabla w \\ &= \sum_{K \in \mathcal{T}} \int_K R_K(u_{\mathcal{T}}) w + \sum_{E \in \mathcal{E}} \int_E R_E(u_{\mathcal{T}}) w. \end{aligned}$$

Together with the Galerkin orthogonality this implies

$$\begin{aligned} & \int_{\Omega} f w + \int_{\Gamma_N} g w - \int_{\Omega} \nabla u_{\mathcal{T}} \cdot \nabla w \\ &= \sum_{K \in \mathcal{T}} \int_K R_K(u_{\mathcal{T}}) (w - w_{\mathcal{T}}) \\ &\quad + \sum_{E \in \mathcal{E}} \int_E R_E(u_{\mathcal{T}}) (w - w_{\mathcal{T}}) \end{aligned}$$

for all  $w \in H_D^1(\Omega)$  and all  $w_{\mathcal{T}} \in S_D^{1,0}(\mathcal{T})$ .

**II.1.7. Upper error bound.** We fix an arbitrary function  $w \in H_D^1(\Omega)$  and choose  $w_{\mathcal{T}} = I_{\mathcal{T}} w$  with the quasi-interpolation operator of Section I.2.11 (p. 21). The Cauchy-Schwarz inequality for integrals and the properties of  $I_{\mathcal{T}}$  then yield

$$\begin{aligned} & \int_{\Omega} f w + \int_{\Gamma_N} g w - \int_{\Omega} \nabla u_{\mathcal{T}} \cdot \nabla w \\ &= \sum_{K \in \mathcal{T}} \int_K R_K(u_{\mathcal{T}}) (w - I_{\mathcal{T}} w) + \sum_{E \in \mathcal{E}} \int_E R_E(u_{\mathcal{T}}) (w - I_{\mathcal{T}} w) \\ &\leq \sum_{K \in \mathcal{T}} \|R_K(u_{\mathcal{T}})\|_K \|w - I_{\mathcal{T}} w\|_K + \sum_{E \in \mathcal{E}} \|R_E(u_{\mathcal{T}})\|_E \|w - I_{\mathcal{T}} w\|_E \\ &\leq \sum_{K \in \mathcal{T}} \|R_K(u_{\mathcal{T}})\|_{K C A_1} h_K \|w\|_{H^1(\bar{\omega}_K)} \\ &\quad + \sum_{E \in \mathcal{E}} \|R_E(u_{\mathcal{T}})\|_{E C A_2} h_E^{\frac{1}{2}} \|w\|_{H^1(\bar{\omega}_E)}. \end{aligned}$$

Invoking the Cauchy-Schwarz inequality for sums this gives

$$\begin{aligned} & \int_{\Omega} f w + \int_{\Gamma_N} g w - \int_{\Omega} \nabla u_{\mathcal{T}} \cdot \nabla w \\ &\leq \max\{C_{A_1}, C_{A_2}\} \left\{ \sum_{K \in \mathcal{T}} h_K^2 \|R_K(u_{\mathcal{T}})\|_K^2 \right\} \end{aligned}$$

$$\left. + \sum_{E \in \mathcal{E}} h_E \|R_E(u_{\mathcal{T}})\|_E^2 \right\}^{\frac{1}{2}} \cdot \left\{ \sum_{K \in \mathcal{T}} \|w\|_{H^1(\bar{\omega}_K)}^2 + \sum_{E \in \mathcal{E}} \|w\|_{H^1(\bar{\omega}_E)}^2 \right\}^{\frac{1}{2}}.$$

In a last step we observe that the shape-regularity of  $\mathcal{T}$  implies

$$\left\{ \sum_{K \in \mathcal{T}} \|w\|_{H^1(\bar{\omega}_K)}^2 + \sum_{E \in \mathcal{E}} \|w\|_{H^1(\bar{\omega}_E)}^2 \right\}^{\frac{1}{2}} \leq c \|w\|_{H^1(\Omega)}$$

with a constant  $c$  which only depends on the shape parameter  $C_{\mathcal{T}}$  of  $\mathcal{T}$  and which takes into account that every element is counted several times on the left-hand side of this inequality.

Combining these estimates with the equivalence of error and residual, we obtain the following upper bound on the error

$$\|u - u_{\mathcal{T}}\|_{H^1(\Omega)} \leq c^* \left\{ \sum_{K \in \mathcal{T}} h_K^2 \|R_K(u_{\mathcal{T}})\|_K^2 + \sum_{E \in \mathcal{E}} h_E \|R_E(u_{\mathcal{T}})\|_E^2 \right\}^{\frac{1}{2}}$$

with

$$c^* = \sqrt{1 + c_{\Omega}^2} \max\{c_{A1}, c_{A2}\}c.$$

The right-hand side of this estimate can be used as an *a posteriori error estimator* since it only involves the known data  $f$  and  $g$ , the solution  $u_{\mathcal{T}}$  of the discrete problem, and the geometrical data of the partition. The above inequality implies that the a posteriori error estimator is *reliable* in the sense that an inequality of the form "error estimator  $\leq$  tolerance" implies that the true error is also less than the tolerance up to the multiplicative constant  $c^*$ . We want to show that the error estimator is also *efficient* in the sense that an inequality of the form "error estimator  $\geq$  tolerance" implies that the true error is also greater than the tolerance possibly up to another multiplicative constant.

For general functions  $f$  and  $g$  the exact evaluation of the integrals occurring on the right-hand side of the above estimate may be prohibitively expensive or even impossible. The integrals then must be approximated by suitable quadrature formulae. Alternatively the functions  $f$  and  $g$  may be approximated by simpler functions, e.g., piecewise polynomial ones, and the resulting integrals be evaluated exactly. Often, both approaches are equivalent.

**II.1.8. Lower error bound.** In order to prove the announced efficiency, we denote for every element  $K$  by  $f_K$  the mean value of  $f$  on  $K$

$$f_K = \frac{1}{|K|} \int_K f dx$$

and for every edge  $E$  on the Neumann boundary by  $g_E$  the mean value of  $g$  on  $E$

$$g_E = \frac{1}{|E|} \int_E g dS.$$

We fix an arbitrary element  $K$  and insert the function

$$w_K = (f_K + \Delta u_{\mathcal{T}})\psi_K$$

in the  $L^2$ -representation of the residual. Taking into account that  $\text{supp } w_K \subset K$  we obtain

$$\int_K R_K(u_{\mathcal{T}})w_K = \int_K \nabla(u - u_{\mathcal{T}}) \cdot \nabla w_K.$$

We add  $\int_K (f_K - f)w_K$  on both sides of this equation and obtain

$$\begin{aligned} \int_K (f_K + \Delta u_{\mathcal{T}})^2 \psi_K &= \int_K (f_K + \Delta u_{\mathcal{T}})w_K \\ &= \int_K \nabla(u - u_{\mathcal{T}}) \cdot \nabla w_K - \int_K (f - f_K)w_K. \end{aligned}$$

The results of Section 1.2.12 (p. 22) imply for the left hand-side of this equation

$$\int_K (f_K + \Delta u_{\mathcal{T}})^2 \psi_K \geq c_{T1}^2 \|f_K + \Delta u_{\mathcal{T}}\|_K^2$$

and for the two terms on its right-hand side

$$\begin{aligned} \int_K \nabla(u - u_{\mathcal{T}}) \cdot \nabla w_K &\leq \|\nabla(u - u_{\mathcal{T}})\|_K \|\nabla w_K\|_K \\ &\leq \|\nabla(u - u_{\mathcal{T}})\|_K c_{T2} h_K^{-1} \|f_K + \Delta u_{\mathcal{T}}\|_K \\ \int_K (f - f_K)w_K &\leq \|f - f_K\|_K \|w_K\|_K \\ &\leq \|f - f_K\|_K \|f_K + \Delta u_{\mathcal{T}}\|_K. \end{aligned}$$

This proves that

$$(II.1.4) \quad h_K \|f_K + \Delta u_{\mathcal{T}}\|_K \leq c_{T1}^{-2} c_{T2} \|\nabla(u - u_{\mathcal{T}})\|_K + c_{T1}^{-2} h_K \|f - f_K\|_K.$$

PDE: Find  $u \in H_0^1(\Omega)$ :  $\int_{\Omega} \partial_x u \cdot \partial_x v = \int_{\Omega} f v \quad \forall v \in H_0^1(\Omega)$

FE equation: Find  $u_h \in H_h \subseteq H_0^1(\Omega)$ :  $\int_{\Omega} \partial_x u_h \cdot \partial_x v_h = \int_{\Omega} f v_h \quad \forall v_h \in H_h$

Poincaré-Friedrich:  $\int_{\Omega} v^2 \leq C \int_{\Omega} (\partial_x v)^2$

$\hookrightarrow$  Corollary:  $\frac{1}{\sqrt{C+1}} \|v\|_{H^1(\Omega)} \leq \sup_{\|w\|_{H^1}=1} \int_{\Omega} \partial_x v \partial_x w \leq \|w\|_{H^1}$

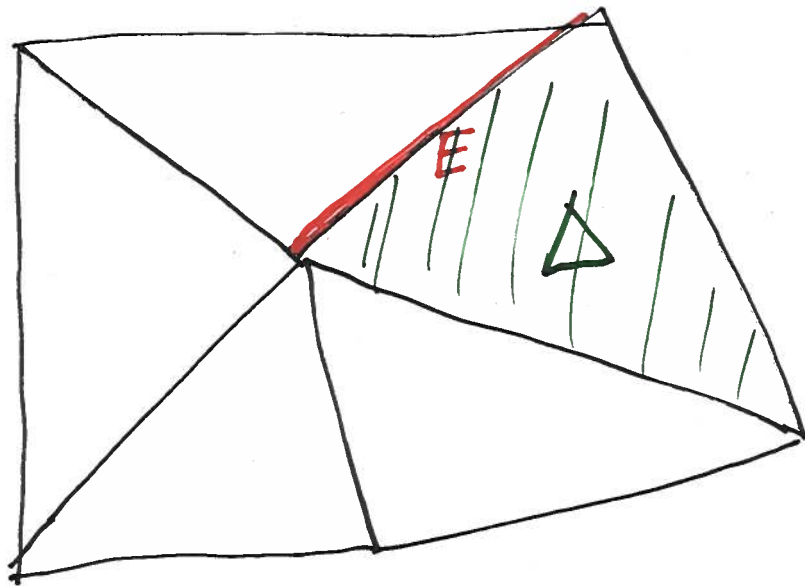
With PDE:  $\int_{\Omega} \partial_x (u - u_h) \partial_x v = \underbrace{\int_{\Omega} f v - \int_{\Omega} \partial_x u_h \partial_x v}_{\text{computable for all } v} =: R(v)$

$\implies \frac{1}{\sqrt{C+1}} \|u - u_h\|_{H^1} \leq \sup_{\|v\|_{H^1}=1} R(v) \leq \|u - u_h\|_{H^1}$

Simplify:  $R(v) \stackrel{\uparrow}{=} \sum_i \int_{x_i}^{x_{i+1}} f(v - v_h) + \sum_i \llbracket \partial_x u_h \rrbracket_{x=x_i} (w(x_i) - w_h(x_i))$   
Galerkin orthogonality

→ Residual error estimate in higher dimension

$$\|u - u_h\|_{H^1(\Omega)} \leq \tilde{C} \left\{ \sum_{\Delta} h_{\Delta}^2 \int_{\Delta} f^2 + \sum_E h_E \int_E \left\| \frac{\partial u_h}{\partial n} \right\|^2 \right\}^{1/2}$$



$\sum_E$  vanishes in 1d