

A POSTERIORI ANALYSIS FOR PDE's.

S. REPIN,

V.A. Steklov Institute of Mathematics in St.-Petersburg

1 Lecture 1. FIRST APPROACHES TO A POSTERIORI ERROR CONTROL

- Error estimation problem in computer simulation
- Mathematical background
- A priori error estimation methods
- Runge's rule
- The estimate of Prager and Synge
- The estimate of Mikhlin
- A posteriori estimates of Ostrowski for contractive mappings
- A posteriori methods based on monotonicity

2 Lecture 2. A POSTERIORI ERROR ESTIMATION METHODS FOR FEM

- Sobolev spaces with negative indices
- Errors and Residuals. First glance
- Residual type estimates for elliptic equations
- Explicit residual method in 1D case
- Explicit residual method in 2D case
- A posteriori error indicators based on post-processing of computed solutions
- General idea
- Post-processing by averaging

- Superconvergence
- Post-processing by equilibration
- A posteriori error estimates constructed with help of adjoint problems

3 FUNCTIONAL A POSTERIORI ESTIMATES FOR A MODEL ELLIPTIC PROBLEM

- Deriving functional a posteriori estimates by the variational method
- Deriving functional a posteriori estimates by the non-variational method
- Properties of functional a posteriori estimates
- How to use the estimates in practice?

4 Two-sided a posteriori estimates for linear elliptic problems

- Problem in the abstract form
- Upper bounds of the error
- Lower bounds of the error
- Computability of two-sided estimates
- Relationships with other methods
- Diffusion problem
- Linear elasticity problem
- Elliptic equations of the 4th order

5 A POSTERIORI ESTIMATES IN NON-ENERGY QUANTITIES

- General principle

- Error estimation in terms of goal-oriented functionals
- A method using post-processing
- A method using functional type two-sided estimates
- Error estimates in terms of seminorms

6 A POSTERIORI ESTIMATES FOR MIXED METHODS

- Mixed formulations of elliptic problems
- Primal Mixed Finite Element Method
- Dual Mixed Finite Element Method
- A priori error estimates for Dual Mixed Finite Element Method
- Functional a posteriori estimates for mixed approximations

7 MIXED FEM ON DISTORTED MESHES

- Distorted meshes
- Extension of fluxes inside a cell
- Well-posedness of discrete problems
- Examples of extension operators
- A priori rate convergence estimates
- Inf-sup condition
- A posteriori estimates

8 A POSTERIORI ESTIMATES FOR PROBLEMS IN THE THEORY OF VISCOUS FLUIDS

- Mathematical models of viscous fluids
- Stokes problem
- Inf–Sup condition
- Existence of a saddle–point
- Estimates of the distance to the set of solenoidal fields
- A posteriori estimates for the Stokes problem
- Estimates for almost incompressible fluids
- Generalizations
- Nonlinear models of viscous fluids
- Antiplane flow of Bingham fluid
- Functional a posteriori estimates for generalized Newtonian fluids

9 A POSTERIORI ESTIMATES FOR NONLINEAR VARIATIONAL PROBLEMS

- Abstract variational problem
- Dual and bidual functionals
- Compound functionals
- Uniformly convex functionals
- General form of the functional a posteriori estimate
- Example. Diffusion problem with Robin boundary conditions

Preface

This is the electronic version of the lecture course

**Mat-5.210 Special Course in Computational Mechanics Autumn 2006:
"A POSTERIORI ANALYSIS FOR PDE's"
<http://math.tkk.fi/teaching/lme/>**

that was prepared for students and PhD students of the Department of Mathematics of Helsinki University of Technology and read in 2006. In general, it is aimed to give an answer to the question "How to verify the accuracy of approximate solutions of partial differential equations computed by various numerical methods?" A posteriori error estimates present a tool able to give an answer to the above question. Nowadays a posteriori estimates form a basis of many powerful numerical techniques and are widely used in computational technologies. Therefore, the purpose of the course is *to discuss the main lines in a posteriori analysis and explain their mathematical foundations.*

First, the "classical" a posteriori error estimation methods developed for finite element approximations are considered. However, the major part of the course is devoted to a new *functional approach* to a posteriori error estimation developed in the last decade. Basic ideas of this approach are first explained on the paradigm of a simple elliptic problem. Further exposition contains applications to particular classes of problems: diffusion, linear elasticity, Stokes, biharmonic, variational inequalities, etc.

The material is based on earlier lectures on a posteriori estimates and adaptive methods that has been read by the author at the University of Houston, USA (2002); University of Jyväskylä, Finland (2003), Radon Institute of Computational and Applied Mathematics (RICAM) in Linz (2005), and at St.-Petersburg Polytechnical University (2000-2003). However, in general, the course is new. A special attention is paid on a posteriori error estimation methods for two important classes of problems that are now in the focus of numerous researches in numerical analysis, namely to *mixed FE approximations* and *approximations in the theory of viscous fluids*.

Full list of references is given at the end of the text, but certain key publications are also cited in the respective places related to the topic discussed.

I wish to express my gratitude to Helsinki University of Technology for the invitation and hospitality. I am grateful to Prof. R. Stenberg for his kind support and interesting discussions that helped to create the course in the present form and also to Mr. A. Niemi and Mrs. Tuula–Donskoi for their help during my visit.

Sergey Repin

Espoo, November 2006

Lecture 1

The goal of the lecture is to
provide a **background** information,
shortly discuss the **a priori error estimation methods**
and to give a concise overview of
first a posteriori error estimation methods

Lecture plan

- Error estimation problem in computer simulation;
- A priori approach to the error analysis for PDE's;
- A posteriori methods developed in 1900–1975:
 - Heuristic Runge's rule;
 - Prager and Synge estimate;
 - Estimate of Mikhlin;
 - Estimates of Ostrowski for contractive mappings;
 - Estimates based on monotonicity (Collatz);

Let us begin with a "philosophic" question:

**WHAT THE NUMBERS COMPUTED
INDEED MEAN?**

To convince yourself that the question stated is worth thinking out, please make

Task 1. "Baby" coupled problem.

$$\begin{aligned} z'' - 9z' - 10z &= 0, & z &= z(x), & x &\in [0, 8], \\ z(0) &= 1, & z'(0) &= \mathbf{a}_{N-1} - \mathbf{a}_N, \end{aligned}$$

where \mathbf{a} is a solution of the system of the dimensionality \mathbf{N}

$$\begin{aligned} \mathbf{B}\mathbf{a} &= \mathbf{f}, & \mathbf{b}_{ij} &= \frac{2\mathbf{S}_i^2\mathbf{S}_j^2}{\pi} \int_0^\pi (\sin(i\xi)\sin(j\xi) + \sin(i+j^2)\xi) d\xi, \\ \mathbf{i}, \mathbf{j} &= 1, 2, \dots, \mathbf{N}, & \mathbf{f}_i &= (\mathbf{i} + 1)^4 \mathbf{i}, & \mathbf{S}_i &= \sum_{k=0}^{+\infty} \left(\frac{\mathbf{i}}{\mathbf{i} + 1} \right)^k. \end{aligned}$$

The task

For $\mathbf{N} = 10, 50, 100, 200$ find $\mathbf{z}(\mathbf{8})$ analytically and compare with numerical results obtained by computing the sums numerically, finding definite integrals with help of quadratures formulas, solving the system of linear simultaneous equations by a numerical method, and integrating the differential equation by a certain (e.g., Euler) method.

I. In the vast majority of cases, exact solutions of differential equations are unknown. We have no other way to use differential equations in the mathematical modeling other than compute approximate solutions and analyze computer simulation results.

II. Approximate solutions contain errors of various nature.

From **I** and **II**, it follows that

III. Verification of the accuracy of approximate solutions a KEY QUESTION.

Errors in mathematical modeling

- ε_1 – error of a mathematical model used
- ε_2 – approximation error arising when a differential model is replaced by a discrete one;
- ε_3 – numerical errors arising when solving a discrete problem.

MODELING ERROR

Let \mathbf{U} be a physical value that characterizes some process and \mathbf{u} be a respective value obtained from the mathematical model. Then the quantity

$$\varepsilon_1 = |\mathbf{U} - \mathbf{u}|$$

is an **error of the mathematical model**.

Mathematical model always presents an "abridged" version of a physical object.

Therefore, $\varepsilon_1 > 0$.

TYPICAL SOURCES OF MODELING ERRORS

- (a) "Second order" phenomena are neglected in a mathematical model.

- (b) Problem data are defined with an uncertainty.

- (c) Dimension reduction is used to simplify a model.

APPROXIMATION ERROR

Let \mathbf{u}_h be a solution on a mesh of the size \mathbf{h} . Then, \mathbf{u}_h encompasses the **approximation error**

$$\epsilon_2 = |\mathbf{u} - \mathbf{u}_h|.$$

Classical error control theory is mainly focused on approximation errors.

NUMERICAL ERRORS

Finite-dimensional problems are also solved approximately, so that instead of u_h we obtain u_h^ε . The quantity

$$\varepsilon_3 = |u_h - u_h^\varepsilon|$$

shows an **error of the numerical algorithm** performed with a concrete computer. This error includes

- roundoff errors,
- errors arising in iteration processes and in numerical integration,
- errors caused by possible defects in computer codes.

Roundoff errors

Numbers in a computer are presented in a **floating point format**:

$$x = \pm \left(\frac{i_1}{q} + \frac{i_2}{q^2} + \dots + \frac{i_k}{q^k} \right) q^\ell, \quad i_s < q.$$

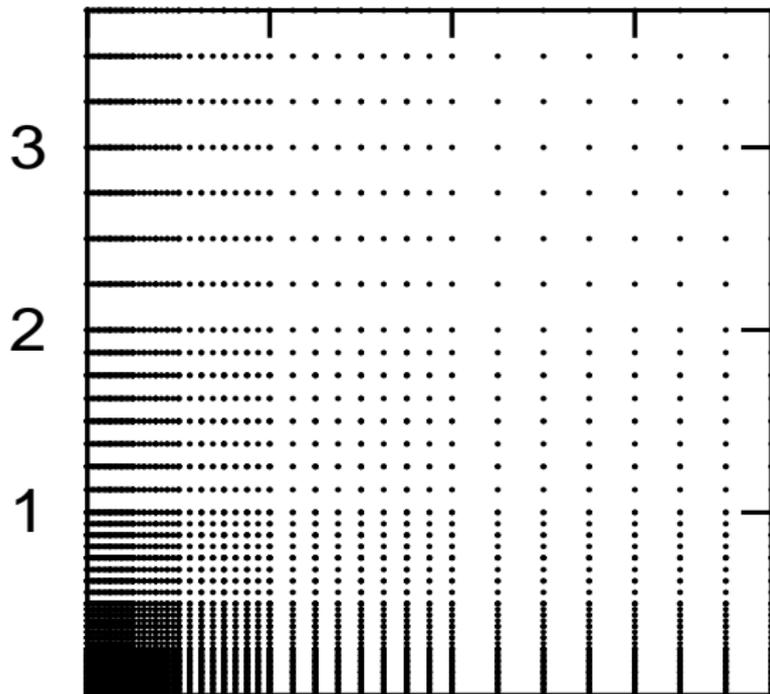
These numbers form the set $R_{q\ell k} \subset \mathbb{R}$.

q is the base of the representation,

$\ell \in [\ell_1, \ell_2]$ is the power.

$R_{q\ell k}$ is not closed with respect to the operations $+$, $-$, $*$!

The set $\mathbf{R}_{q\ell k} \times \mathbf{R}_{q\ell k}$



Example

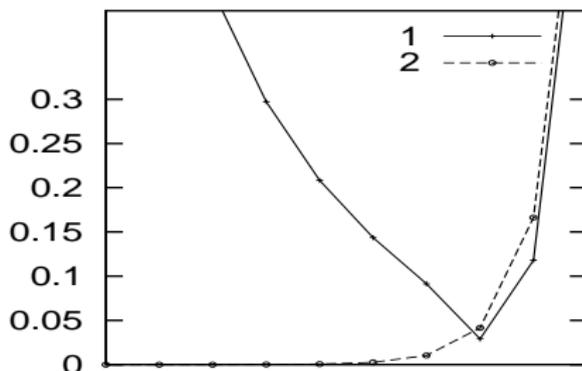
$$\mathbf{k} = 3, \quad \mathbf{a} = \left(\frac{1}{2} + 0 + 0\right) * 2^5, \quad \mathbf{b} = \left(\frac{1}{2} + 0 + 0\right) * 2^1$$
$$\mathbf{b} \Rightarrow \left(0 + \frac{1}{2} + 0\right) * 2^2 \Rightarrow \left(0 + 0 + \frac{1}{2}\right) * 2^3 \Rightarrow (0 + 0 + 0) * 2^4$$

$$\mathbf{a} + \mathbf{b} = \mathbf{a}!!!$$

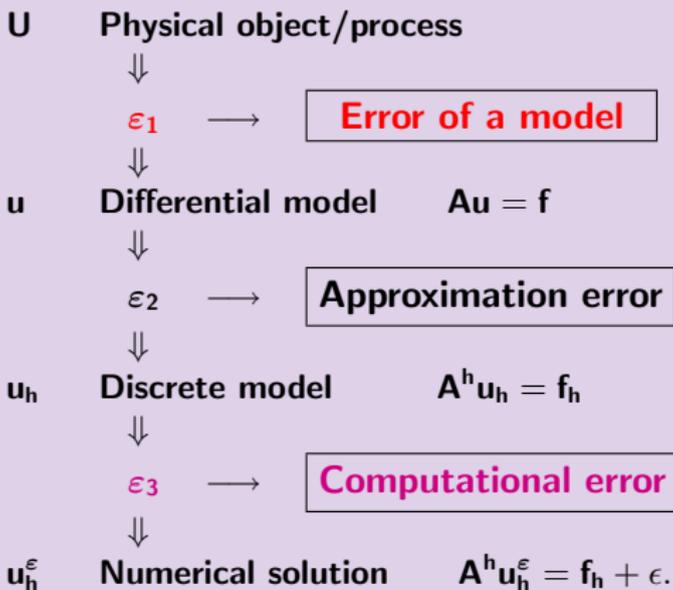
Definition. The smallest floating point number which being added to 1 gives a quantity other than 1 is called **the machine accuracy**.

Numerical integration

$$\int_b^a f(x) dx \cong \sum_{i=1}^n c_i f(x_i) h = \sum_{i=1}^{n/2} c_i^{\sim 1} f(x_i) h + c_{n/2+1}^{\sim \delta} f(x_{n/2+1}) h + \dots$$



Errors in computer simulation



Two principal relations

I. Computations on the basis of a reliable (certified) model. Here ϵ_1 is assumed to be small and u_h^ϵ gives a desired information on U .

$$\|U - u_h^\epsilon\| \leq \epsilon_1 + \boxed{\epsilon_2 + \epsilon_3}. \quad (1)$$

II. Verification of a mathematical model. Here physical data U and numerical data u_h^ϵ are compared to judge on the quality of a mathematical model

$$\|\epsilon_1\| \leq \|U - u_h^\epsilon\| + \boxed{\epsilon_2 + \epsilon_3}. \quad (2)$$

Thus, two major problems of mathematical modeling, namely,

- reliable computer simulation,
- verification of mathematical models by comparing physical and mathematical experiments,

require efficient methods able to provide
COMPUTABLE AND REALISTIC
estimates of $\epsilon_2 + \epsilon_3$.

What is u and what is $\|\cdot\|$?

If we start a more precise investigation, then it is necessary to answer the question

What is a solution to a boundary-value problem?

Example.

$$\frac{\partial^2 \mathbf{u}}{\partial \mathbf{x}_1^2} + \frac{\partial^2 \mathbf{u}}{\partial \mathbf{x}_2^2} + \mathbf{f} = \mathbf{0}, \quad \mathbf{u} = \mathbf{u}_0 \text{ on } \partial\Omega \quad \boxed{\exists \mathbf{u}?$$

It is not a trivial question, so that about one hundred years passed before mathematicians have found an appropriate concept for PDE's.

Without proper understanding of a mathematical model no real modeling can be performed. Indeed,

If we are not sure that a solution u exists then what we try to approximate numerically?

If we do not know to which class of functions u belongs to, then we cannot properly define the measure for the accuracy of computed approximations.

Thus, we need to recall a

CONCISE MATHEMATICAL BACKGROUND

Vectors and tensors

\mathbb{R}^n contains real n -vectors. $\mathbb{M}^{n \times m}$ contains $n \times m$ matrices and $\mathbb{M}_s^{n \times n}$ contains $n \times n$ symmetric matrices (tensors) with real entries.

$$\mathbf{a} \cdot \mathbf{b} = \sum_{i=1}^n \mathbf{a}_i \mathbf{b}_i \in \mathbb{R}, \quad \mathbf{a}, \mathbf{b} \in \mathbb{R}^n \quad (\text{scalar product of vectors}),$$

$$\mathbf{a} \otimes \mathbf{b} = \{\mathbf{a}_i \mathbf{b}_j\} \in \mathbb{M}^{n \times n} \quad (\text{tensor product of vectors}),$$

$$\boldsymbol{\sigma} : \boldsymbol{\varepsilon} = \sum_{i,j=1}^n \boldsymbol{\sigma}_{ij} \boldsymbol{\varepsilon}_{ij} \in \mathbb{R}, \quad \boldsymbol{\sigma}, \boldsymbol{\varepsilon} \in \mathbb{M}^{n \times n} \quad (\text{scalar product of tensors}).$$

$$|\mathbf{a}| := \sqrt{\mathbf{a} \cdot \mathbf{a}}, \quad |\boldsymbol{\sigma}| := \sqrt{\boldsymbol{\sigma} : \boldsymbol{\sigma}},$$

Unit matrix is denoted by \mathbb{I} . If $\boldsymbol{\tau} \in \mathbb{M}^{n \times n}$, then $\boldsymbol{\tau}^D = \boldsymbol{\tau} - \frac{1}{n} \mathbb{I}$ is the deviator of $\boldsymbol{\tau}$.

Spaces of functions

Let Ω be an open bounded domain in \mathbb{R}^n with Lipschitz continuous boundary.

$\mathbf{C}^k(\Omega)$ – k times continuously differentiable functions.

$\mathbf{C}_0^k(\Omega)$ – k times continuously differentiable functions vanishing at the boundary $\partial\Omega$.

$\mathbf{C}_0^\infty(\Omega)$ – k smooth functions with compact supports in Ω .

$\mathbf{L}^p(\Omega)$ – summable functions with finite norm

$$\|\mathbf{g}\|_{p,\Omega} = \|\mathbf{g}\|_p = \left(\int_{\Omega} |\mathbf{g}|^p \right)^{1/p}.$$

For $\mathbf{L}^2(\Omega)$ the norm is denoted by $\|\cdot\|$.

If \mathbf{g} is a vector (tensor)-valued function, then the respective spaces are denoted by

$\mathbf{C}^k(\Omega, \mathbb{R}^n)$ ($\mathbf{C}^k(\Omega, \mathbb{M}^{n \times n})$),

$\mathbf{L}^p(\Omega, \mathbb{R}^n)$ ($\mathbf{L}^p(\Omega, \mathbb{M}^{n \times n})$)

with similar norms.

We say that \mathbf{g} is **locally integrable** in Ω and write $\mathbf{f} \in \mathbf{L}^{1,loc}(\Omega)$, if $\mathbf{g} \in \mathbf{L}^1(\omega)$ for any $\omega \subset\subset \Omega$. Similarly, one can define the space $\mathbf{L}^{p,loc}(\Omega)$ that consists of functions locally integrable with degree $\mathbf{p} \geq \mathbf{1}$.

Generalized derivatives

Let $\mathbf{f}, \mathbf{g} \in \mathbf{L}^{1,\text{loc}}(\Omega)$ and

$$\int_{\Omega} \mathbf{g} \varphi \, d\mathbf{x} = - \int_{\Omega} \mathbf{f} \frac{\partial \varphi}{\partial x_i} \, d\mathbf{x}, \quad \forall \varphi \in \mathring{\mathbf{C}}^1(\Omega).$$

Then \mathbf{g} is called a **generalized derivative** (in the sense of Sobolev) of \mathbf{f} with respect to x_i and we write

$$\mathbf{g} = \frac{\partial \mathbf{f}}{\partial x_i}.$$

Higher order generalized derivatives

If $\mathbf{f}, \mathbf{g} \in \mathbf{L}^{1,\text{loc}}(\Omega)$ and

$$\int_{\Omega} \mathbf{g} \varphi \, d\mathbf{x} = \int_{\Omega} \mathbf{f} \frac{\partial^2 \varphi}{\partial x_i \partial x_j} \, d\mathbf{x}, \quad \forall \varphi \in \mathring{\mathbf{C}}^2(\Omega),$$

then \mathbf{g} is a generalized derivative of \mathbf{f} with respect to x_i and x_j . For generalized derivatives we keep the classical notation and write

$$\mathbf{g} = \partial^2 \mathbf{f} / \partial x_i \partial x_j = \mathbf{f}_{,ij}.$$

If \mathbf{f} is differentiable in the classical sense, then its generalized derivatives coincide with the classical ones !

To extend this definition further, we use the multi-index notation and write $\mathbf{D}^\alpha \mathbf{f}$ in place of $\partial^k \mathbf{f} / \partial x_1^{\alpha_1} \partial x_2^{\alpha_2} \dots \partial x_n^{\alpha_n}$.

Definition

Let $\mathbf{f}, \mathbf{g} \in \mathbf{L}^{1, \text{loc}}(\Omega)$ and

$$\int_{\Omega} \mathbf{g} \varphi \, d\mathbf{x} = (-1)^{|\alpha|} \int_{\Omega} \mathbf{f} \mathbf{D}^\alpha \varphi \, d\mathbf{x}, \quad \forall \varphi \in \mathring{\mathbf{C}}^k(\Omega).$$

Then, \mathbf{g} is called a **generalized derivative** of \mathbf{f} of degree

$$|\alpha| := \alpha_1 + \alpha_2 + \dots + \alpha_n$$

and we write

$$\mathbf{g} = \mathbf{D}^\alpha \mathbf{f}.$$

Sobolev spaces

S. L. Sobolev. *Some Applications of Functional Analysis in Mathematical Physics*, Izdat. Leningrad. Gos. Univ., Leningrad, 1955, English version: Translation of Mathematical Monographs, Volume 90, American Mathematical Society, Providence, RI, 1991.

Definition

The spaces of functions that have integrable generalized derivatives up to a certain order are called **Sobolev spaces**. $\mathbf{f} \in \mathbf{W}^{1,p}(\Omega)$ if $\mathbf{f} \in \mathbf{L}^p$ and all the generalized derivatives of \mathbf{f} of the first order belong to L^p , i.e.,

$$\mathbf{f}_{,i} = \frac{\partial \mathbf{f}}{\partial x_i} \in \mathbf{L}^p(\Omega).$$

The norm in $\mathbf{W}^{1,p}$ is defined as follows:

$$\|\mathbf{f}\|_{1,p,\Omega} := \left(\int_{\Omega} (|\mathbf{f}|^p + \sum_{i=1}^n |\mathbf{f}_{,i}|^p) \mathbf{d}\mathbf{x} \right)^{1/p}.$$

All the other Sobolev spaces are defined quite similarly: $\mathbf{f} \in \mathbf{W}^{k,p}(\Omega)$ if all generalized derivatives up to the order k are integrable with power p and the quantity

$$\|\mathbf{f}\|_{k,p,\Omega} := \left(\int_{\Omega} \sum_{|\alpha| \leq k} |\mathbf{D}^{\alpha} \mathbf{f}|^p \, d\mathbf{x} \right)^{1/p}$$

is finite. For the Sobolev spaces $\mathbf{W}^{k,2}(\Omega)$ we also use a simplified notation $\mathbf{H}^k(\Omega)$.

Sobolev spaces of vector- and tensor-valued functions are introduced by obvious extensions of the above definitions. We denote them by $\mathbf{W}^{k,p}(\Omega, \mathbb{R}^n)$ and $\mathbf{W}^{k,p}(\Omega, \mathbb{M}^{n \times n})$, respectively.

Embedding Theorems

Relationships between the Sobolev spaces and $\mathbf{L}^p(\Omega)$ and $\mathbf{C}^k(\Omega)$ are given by **Embedding Theorems**.

If $p, q \geq 1$, $\ell > 0$ and $\ell + \frac{n}{q} \geq \frac{n}{p}$, then $\mathbf{W}^{\ell,p}(\Omega)$ is continuously embedded in $\mathbf{L}^q(\Omega)$. Moreover, if $\ell + \frac{n}{q} > \frac{n}{p}$, then the embedding operator is compact.

If $\ell - k > \frac{n}{p}$, then $\mathbf{W}^{\ell,p}(\Omega)$ is compactly embedded in $\mathbf{C}^k(\overline{\Omega})$.

Traces

The functions in Sobolev spaces have counterparts on $\partial\Omega$ called **traces**. Thus, there exist some bounded operators mapping the functions defined in Ω to functions defined on the boundary, e.g.,

$$\gamma : \mathbf{H}^1(\Omega) \rightarrow \mathbf{L}^2(\partial\Omega)$$

is called the **trace operator** if it satisfies the following conditions:

$$\begin{aligned}\gamma\mathbf{v} &= \mathbf{v} |_{\partial\Omega}, & \forall \mathbf{v} \in \mathbf{C}^1(\Omega), \\ \|\gamma\mathbf{v}\|_{2,\partial\Omega} &\leq \mathbf{c}\|\mathbf{v}\|_{1,2,\Omega},\end{aligned}$$

where \mathbf{c} is a positive constant independent of \mathbf{v} . From these relations, we observe that such a trace is a natural generalization of the trace defined for a continuous function.

It was established that $\gamma\mathbf{v}$ forms a subset of $\mathbf{L}^2(\partial\Omega)$, which is the space $\mathbf{H}^{1/2}(\partial\Omega)$. The functions from other Sobolev spaces also are known to have traces in Sobolev spaces with fractional indices.

Henceforth, we understand the boundary values of functions in the sense of traces, so that

$$\mathbf{u} = \psi \quad \text{on } \partial\Omega$$

means that the trace $\gamma\mathbf{u}$ of a function \mathbf{u} defined in Ω coincides with a given function ψ defined on $\partial\Omega$.

All the spaces of functions that have zero traces on the boundary are marked by the symbol \circ (e.g., $\overset{\circ}{\mathbf{W}}^{1,p}(\Omega)$ and $\overset{\circ}{\mathbf{H}}^1(\Omega)$).

Inequalities

1. Friederichs-Steklov inequality.

$$\|\mathbf{w}\| \leq \mathbf{C}_\Omega \|\nabla \mathbf{w}\|, \quad \forall \mathbf{w} \in \mathring{\mathbf{H}}^1(\Omega), \quad (3)$$

2. Poincaré inequality.

$$\|\mathbf{w}\| \leq \tilde{\mathbf{C}}_\Omega \|\nabla \mathbf{w}\|, \quad \forall \mathbf{w} \in \tilde{\mathbf{H}}^1(\Omega), \quad (4)$$

where $\tilde{\mathbf{H}}^1(\Omega)$ is a subset of H^1 of functions with zero mean.

3. Korn's inequality.

$$\int_{\Omega} (|\mathbf{v}|^2 + |\varepsilon(\mathbf{v})|^2) \, d\mathbf{x} \geq \mu_\Omega \|\mathbf{v}\|_{1,2,\Omega}^2, \quad \forall \mathbf{v} \in \mathbf{H}^1(\Omega, \mathbb{R}^n), \quad (5)$$

Generalized solutions

The concept of generalized solutions to PDE's came from

Petrov-Bubnov-Galerkin method.

B. G. Galerkin. Beams and plates. Series in some questions of elastic equilibrium of beams and plates (approximate translation of the title from Russian). *Vestnik Ingenerov, St.-Peterburg*, 19(1915), 897-908.

$$\int_{\Omega} (\Delta \mathbf{u} + \mathbf{f}) \mathbf{w} \, dx = 0 \quad \forall \mathbf{w}$$

Integration by parts leads to the so-called **generalized formulation** of the problem: find $\mathbf{u} \in \mathring{\mathbf{H}}^1(\Omega) + \mathbf{u}_0$ such that

$$\int_{\Omega} \nabla \mathbf{u} \cdot \nabla \mathbf{w} \, dx = \int_{\Omega} \mathbf{f} \mathbf{w} \, dx \quad \forall \mathbf{w} \in \mathring{\mathbf{H}}^1(\Omega)$$

This idea admits wide extensions to many differential equations, see e.g.,
[O. A. Ladyzhenskaya, *The boundary value problems of mathematical physics*. Springer-Verlag, New York, 1985](#)

Definition

A symmetric form $\mathbf{B} : \mathbf{V} \times \mathbf{V} \rightarrow \mathbf{R}$, where V is a Hilbert space, called *V - elliptic* if $\exists c_1 > 0, c_2 > 0$ such that

$$\mathbf{B}(\mathbf{u}, \mathbf{u}) \geq c_1 \|\mathbf{u}\|^2, \quad \forall \mathbf{u} \in \mathbf{V}$$

$$|\mathbf{B}(\mathbf{u}, \mathbf{v})| \leq c_2 \|\mathbf{u}\| \|\mathbf{v}\|, \quad \forall \mathbf{u}, \mathbf{v} \in \mathbf{V}$$

General formulation for linear PDE's is: for a certain linear continuous functional \mathbf{f} (from the space \mathbf{V}^* topologically dual to \mathbf{V}) find \mathbf{u} such that

$$\mathbf{B}(\mathbf{u}, \mathbf{w}) = \langle \mathbf{f}, \mathbf{w} \rangle \quad \mathbf{w} \in \mathbf{V}.$$

Existence of a solution

Usually, existence is proved by

Lax-Milgram Lemma For a bilinear form \mathbf{B} there exists a linear bounded operator $\mathbf{A} \in \mathcal{L}(\mathcal{V}, \mathcal{V})$ such that

$$\mathbf{B}(\mathbf{u}, \mathbf{v}) = (\mathbf{A}\mathbf{u}, \mathbf{v}), \quad \forall \mathbf{u}, \mathbf{v} \in \mathbf{V}$$

It has an inverse $\mathbf{A}^{-1} \in \mathcal{L}(\mathcal{V}, \mathcal{V})$, such that $\|\mathbf{A}\| \leq \mathbf{c}_2$, $\|\mathbf{A}^{-1}\| \leq \frac{1}{\mathbf{c}_1}$.

We will follow another *modus operandi* !!!.

Variational approach

Lemma

If $J : K \rightarrow \mathbf{R}$ is convex, continuous and coercive, i.e.,

$$J(\mathbf{w}) \rightarrow +\infty \quad \text{as } \|\mathbf{w}\|_{\mathbf{V}} \rightarrow +\infty$$

and K is a convex closed subset of a reflexive space \mathbf{V} , then the problem

$$\inf_{\mathbf{w} \in K} J(\mathbf{w})$$

has a **minimizer** \mathbf{u} . If J is strictly convex, then the minimizer is **unique**.

See, e.g., I. Ekeland and R. Temam. *Convex analysis and variational problems*. North-Holland, Amsterdam, 1976.

Coercivity

Take $J(\mathbf{w}) = \frac{1}{2}B(\mathbf{w}, \mathbf{w}) - \langle \mathbf{f}, \mathbf{w} \rangle$ and let K be a certain subspace. Then

$$\frac{1}{2}B(\mathbf{w}, \mathbf{w}) \geq c_1 \|\mathbf{w}\|_V^2, \quad |\langle \mathbf{f}, \mathbf{w} \rangle| \leq \|\mathbf{f}\|_{V^*} \|\mathbf{w}\|_V.$$

We see, that

$$J(\mathbf{w}) \geq c_1 \|\mathbf{w}\|_V^2 - \|\mathbf{f}\|_{V^*} \|\mathbf{w}\|_V \rightarrow +\infty \quad \text{as } \|\mathbf{w}\|_V \rightarrow +\infty$$

Since J is strictly convex and continuous we conclude that a minimizer exists and unique.

Useful algebraic relation

First we present the algebraic identity

$$\begin{aligned}
 \frac{1}{2}\mathbf{B}(\mathbf{u} - \mathbf{v}, \mathbf{u} - \mathbf{v}) &= \frac{1}{2}\mathbf{B}(\mathbf{v}, \mathbf{v}) - \langle \mathbf{f}, \mathbf{v} \rangle + \\
 &+ \langle \mathbf{f}, \mathbf{u} \rangle - \frac{1}{2}\mathbf{B}(\mathbf{u}, \mathbf{u}) - \mathbf{B}(\mathbf{u}, \mathbf{v} - \mathbf{u}) + \langle \mathbf{f}, \mathbf{v} - \mathbf{u} \rangle = \\
 &= \mathbf{J}(\mathbf{v}) - \mathbf{J}(\mathbf{u}) - \mathbf{B}(\mathbf{u}, \mathbf{v} - \mathbf{u}) + \langle \mathbf{f}, \mathbf{v} - \mathbf{u} \rangle
 \end{aligned} \tag{6}$$

From this identity we derive two important results:

- (a) Minimizer \mathbf{u} satisfies $\mathbf{B}(\mathbf{u}, \mathbf{w}) = \langle \mathbf{f}, \mathbf{w} \rangle \forall \mathbf{w}$;
- (b) Error is subject to the difference of functionals.

Let us show (a), i.e., that from (6) it follows the identity

$$\mathbf{B}(\mathbf{u}, \mathbf{v} - \mathbf{u}) = \langle \mathbf{f}, \mathbf{v} - \mathbf{u} \rangle \quad \forall \mathbf{v} \in \mathbf{K},$$

which is $\mathbf{B}(\mathbf{u}, \mathbf{w}) = \langle \mathbf{f}, \mathbf{w} \rangle$ if set $\mathbf{w} = \mathbf{v} - \mathbf{u}$. Indeed, assume the opposite, i.e. $\exists \tilde{\mathbf{v}} \in \mathbf{K}$ such that

$$\mathbf{B}(\mathbf{u}, \tilde{\mathbf{v}} - \mathbf{u}) - \langle \mathbf{f}, \tilde{\mathbf{v}} - \mathbf{u} \rangle = \delta > 0 \quad (\tilde{\mathbf{v}} \neq \mathbf{u}!)$$

Set $\tilde{\mathbf{v}} := \mathbf{u} + \alpha(\tilde{\mathbf{v}} - \mathbf{u})$, $\alpha \in \mathbb{R}$. Then $\tilde{\mathbf{v}} - \mathbf{u} = \alpha(\tilde{\mathbf{v}} - \mathbf{u})$ and

$$\begin{aligned} \frac{1}{2} \mathbf{B}(\mathbf{u} - \tilde{\mathbf{v}}, \mathbf{u} - \tilde{\mathbf{v}}) + \mathbf{B}(\mathbf{u}, \tilde{\mathbf{v}} - \mathbf{u}) - \langle \mathbf{f}, \tilde{\mathbf{v}} - \mathbf{u} \rangle &= \\ &= \frac{\alpha^2}{2} \mathbf{B}(\tilde{\mathbf{v}} - \mathbf{u}, \tilde{\mathbf{v}} - \mathbf{u}) + \alpha \delta = \mathbf{J}(\tilde{\mathbf{v}}) - \mathbf{J}(\mathbf{u}) \geq 0 \end{aligned}$$

However, for arbitrary α such an inequality cannot be true. Denote $\mathbf{a} = \mathbf{B}(\tilde{\mathbf{v}} - \mathbf{u}, \tilde{\mathbf{v}} - \mathbf{u})$. Then in the left-hand side we have a function $\frac{1}{2}\alpha^2 \mathbf{a}^2 + \alpha \delta$, which always attains negative values for certain α . For example, set $\alpha = -\delta/\mathbf{a}^2$. Then, the left-hand side is equal to $-\frac{1}{2}\delta^2/\mathbf{a}^2 < 0$ and we arrive at a contradiction.

A priori approach to the error control problem

Error estimate

Now, we show (b). From

$$\begin{aligned}\frac{1}{2}\mathbf{B}(\mathbf{u} - \mathbf{v}, \mathbf{u} - \mathbf{v}) &= \\ &= \mathbf{J}(\mathbf{v}) - \mathbf{J}(\mathbf{u}) - \mathbf{B}(\mathbf{u}, \mathbf{v} - \mathbf{u}) + \langle \mathbf{f}, \mathbf{v} - \mathbf{u} \rangle\end{aligned}$$

we obtain the error estimate:

$$\frac{1}{2}\mathbf{B}(\mathbf{u} - \mathbf{v}, \mathbf{u} - \mathbf{v}) = \mathbf{J}(\mathbf{v}) - \mathbf{J}(\mathbf{u}). \quad (7)$$

See **S. G. Mikhlin**. *Variational methods in mathematical physics*. Pergamon, Oxford, 1964.

which immediately gives the **projection estimate**

Projection estimate

Let \mathbf{u}_h be a minimizer of \mathbf{J} on $\mathbf{K}_h \subset \mathbf{K}$. Then

$$\begin{aligned}\frac{1}{2}\mathbf{B}(\mathbf{u} - \mathbf{u}_h, \mathbf{u} - \mathbf{u}_h) &= \mathbf{J}(\mathbf{u}_h) - \mathbf{J}(\mathbf{u}) \leq \mathbf{J}(\mathbf{v}_h) - \mathbf{J}(\mathbf{u}) = \\ &= \frac{1}{2}\mathbf{B}(\mathbf{u} - \mathbf{v}_h, \mathbf{u} - \mathbf{v}_h) \quad \forall \mathbf{v}_h \in \mathbf{K}_h.\end{aligned}$$

and we observe that

$$\boxed{\mathbf{B}(\mathbf{u} - \mathbf{u}_h, \mathbf{u} - \mathbf{u}_h) = \inf_{\mathbf{v}_h \in \mathbf{K}_h} \mathbf{B}(\mathbf{u} - \mathbf{v}_h, \mathbf{u} - \mathbf{v}_h)} \quad (8)$$

Projection type estimates serve a basis for deriving **a priori convergence estimates**.

Interpolation in Sobolev spaces

Two key points: PROJECTION ESTIMATE and INTERPOLATION IN SOBOLEV SPACES.

Interpolation theory investigates the difference between a function in a Sobolev space and its piecewise polynomial interpolant. Basic estimate on a simplex \mathbf{T}_h is

$$|\mathbf{v} - \Pi_h \mathbf{v}|_{m,t,\mathbf{T}_h} \leq \mathbf{C}(m, n, t) \left(\frac{h}{\rho}\right)^m h^{2-m} \|\mathbf{v}\|_{2,t,\mathbf{T}_h},$$

and on the whole domain

$$|\mathbf{v} - \Pi_h \mathbf{v}|_{m,t,\Omega_h} \leq \mathbf{C} h^{2-m} \|\mathbf{v}\|_{2,t,\Omega_h}.$$

Here h is the element size and ρ is the inscribed ball diameter.

Asymptotic convergence estimates

Typical case is $\mathbf{m} = \mathbf{1}$ and $\mathbf{t} = \mathbf{2}$. Since

$$\mathbf{B}(\mathbf{u} - \mathbf{u}_h, \mathbf{u} - \mathbf{u}_h) \leq \mathbf{B}(\mathbf{u} - \mathbf{\Pi}_h \mathbf{u}, \mathbf{u} - \mathbf{\Pi}_h \mathbf{u}) \leq \mathbf{c}_2 \|\mathbf{u} - \mathbf{\Pi}_h \mathbf{u}\|^2$$

for

$$\mathbf{B}(\mathbf{w}, \mathbf{w}) = \int_{\Omega} \nabla \mathbf{w} \cdot \nabla \mathbf{w} \, dx$$

we find that

$$\|\nabla(\mathbf{u} - \mathbf{u}_h)\| \leq \mathbf{C}h |\mathbf{u}|_{2,2,\Omega}.$$

provided that

- Exact solution is H^2 – regular;
- \mathbf{u}_h is the Galerkin approximation;
- Elements do not "degenerate" in the refinement process.

A priori convergence estimates cannot guarantee that the error **monotonically decreases** as $h \rightarrow 0$.

Besides, in practice we are interested in the error of a **concrete approximation on a particular mesh**. Asymptotic estimates could hardly be helpful in such a context because, in general, the constant **C** serves for the *whole class* of approximate solutions of a particular type. Typically it is either unknown or highly overestimated.

A priori convergence estimates have mainly a theoretical value: they show that an approximation method is correct "in principle.

For these reasons, a quite different approach to error control is rapidly developing. Nowadays it has already formed a new direction:

A Posteriori Error Control for PDE's

A posteriori error estimation methods developed in 1900-1975

Runge's rule

At the end of 19th century a heuristic error control method was suggested by C. Runge who investigated numerical integration methods for ordinary differential equations.

Heuristic rule of C. Runge

If the difference between two approximate solutions computed on a coarse mesh \mathcal{T}_h with mesh size \mathbf{h} and refined mesh $\mathcal{T}_{h_{ref}}$ with mesh size \mathbf{h}_{ref} (e.g., $\mathbf{h}_{ref} = \mathbf{h}/2$) has become small, then both $\mathbf{u}_{h_{ref}}$ and \mathbf{u}_h are probably close to the exact solution.

In other words, this rule can be formulated as follows:

If $[\mathbf{u}_h - \mathbf{u}_{h_{ref}}]$ is small then $\mathbf{u}_{h_{ref}}$ is close to \mathbf{u}

where $[\cdot]$ is a certain functional or mesh-dependent norm.

Also, the quantity $[\mathbf{u}_h - \mathbf{u}_{h_{ref}}]$ can be viewed (in terms of modern terminology) as a certain **a posteriori error indicator**.

Runge's heuristic rule is simple and was easily accepted by numerical analysts.

However, if we do not properly define the quantity $[\cdot]$, for which $[\mathbf{u}_h - \mathbf{u}_{h_{\text{ref}}}]$ is small, then the such a principle may be not true.

One can present numerous examples where **two subsequent elements of an approximation sequence are close to each other, but far from a certain joint limit.** For example, such cases often arise in the minimization (maximization) of functionals with "saturation" type behavior or with a "sharp-well" structure. Also, the rule may lead to a wrong presentation if, e.g., the refinement has not been properly done, so that new trial functions were added only in subdomains where an approximation is almost coincide with the true solution. Then two subsequent approximations may be very close, but at the same time not close to the exact solution.

Also, in practice, we need to know precisely what the word "close" means, i.e. we need to have a more concrete presentation on the error. For example, it would be useful to establish the following rule:

$$\text{If } [\mathbf{u}_h - \mathbf{u}_{h,\text{ref}}] \leq \varepsilon \text{ then } \|\mathbf{u}_h - \mathbf{u}\| \leq \delta(\varepsilon),$$

where the function $\delta(\varepsilon)$ is known and computable.

In subsequent lectures we will see that for a wide class of boundary-value problems it is indeed possible to derive such type generalizations of the Runge's rule.

Prager and Synge estimates

W. Prager and J. L. Synge. Approximation in elasticity based on the concept of function spaces, *Quart. Appl. Math.* 5(1947)

Prager and Synge derived an estimate on the basis of purely geometrical grounds. In modern terms, there result for the problem

$$\begin{aligned}\Delta \mathbf{u} + \mathbf{f} &= \mathbf{0}, & \text{in } \Omega, \\ \mathbf{u} &= \mathbf{0}, & \text{on } \partial\Omega\end{aligned}$$

reads as follows:

$$\|\nabla(\mathbf{u} - \mathbf{v})\|^2 + \|\nabla\mathbf{u} - \boldsymbol{\tau}\|^2 = \|\nabla\mathbf{v} - \boldsymbol{\tau}\|^2,$$

where $\boldsymbol{\tau}$ is a function satisfying the equation $\operatorname{div}\boldsymbol{\tau} + \mathbf{f} = \mathbf{0}$. We can easily prove it by the **orthogonality relation**

$$\int_{\Omega} \nabla(\mathbf{u} - \mathbf{v}) \cdot (\nabla\mathbf{u} - \boldsymbol{\tau}) \, d\mathbf{x} = 0 \quad (\operatorname{div}(\nabla\mathbf{u} - \boldsymbol{\tau}) = \mathbf{0}!).$$

Estimate of Mikhlin

S. G. Mikhlin. *Variational methods in mathematical physics*. Pergamon, Oxford, 1964.

A similar estimate was derived by **variational arguments** (see Lecture 1). It is as follows:

$$\frac{1}{2} \|\nabla(\mathbf{u} - \mathbf{v})\|^2 \leq \mathbf{J}(\mathbf{v}) - \mathbf{infJ},$$

where

$$\mathbf{J}(\mathbf{v}) := \frac{1}{2} \|\nabla \mathbf{v}\|^2 - (\mathbf{f}, \mathbf{v}), \quad \mathbf{infJ} := \inf_{\mathbf{v} \in \overset{\circ}{\mathbf{H}}_1(\Omega)} \mathbf{J}(\mathbf{v}).$$

Dual problem

Since

$$\inf \mathbf{J} = \sup_{\boldsymbol{\tau} \in \mathbf{Q}_f} \left\{ -\frac{1}{2} \|\boldsymbol{\tau}\|^2 \right\},$$

where

$$\mathbf{Q}_f := \left\{ \boldsymbol{\tau} \in \mathbf{L}_2(\Omega, \mathbf{R}^d) \mid \int_{\Omega} \boldsymbol{\tau} \cdot \nabla \mathbf{w} \, dx = \int_{\Omega} f \mathbf{w} \, dx \quad \forall \mathbf{w} \in \mathring{\mathbf{H}}^1 \right\},$$

we find that

$$\frac{1}{2} \|\nabla(\mathbf{u} - \mathbf{v})\|^2 \leq \mathbf{J}(\mathbf{v}) + \frac{1}{2} \|\boldsymbol{\tau}\|^2, \quad \forall \boldsymbol{\tau} \in \mathbf{Q}_f.$$

Since

$$\begin{aligned} \mathbf{J}(\mathbf{v}) + \frac{1}{2}\|\boldsymbol{\tau}\|^2 &= \frac{1}{2}\|\nabla\mathbf{v}\|^2 - \int_{\Omega} \mathbf{f}\mathbf{v} \, \mathbf{d}\mathbf{x} + \frac{1}{2}\|\boldsymbol{\tau}\|^2 = \\ &= \frac{1}{2}\|\nabla\mathbf{v}\|^2 - \int_{\Omega} \boldsymbol{\tau} \cdot \nabla\mathbf{v} \, \mathbf{d}\mathbf{x} + \frac{1}{2}\|\boldsymbol{\tau}\|^2 = \\ &= \frac{1}{2}\|\nabla\mathbf{v} - \boldsymbol{\tau}\|^2 \end{aligned}$$

we arrive at the estimate

$$\frac{1}{2}\|\nabla(\mathbf{u} - \mathbf{v})\|^2 \leq \frac{1}{2}\|\nabla\mathbf{v} - \boldsymbol{\tau}\|^2, \quad \forall \boldsymbol{\tau} \in \mathbf{Q}_f. \quad (9)$$

Difficulties

Estimates of Prager and Synge and of Mikhlin are valid for any $\mathbf{v} \in \mathring{\mathbf{H}}_1(\Omega)$, so that, formally, that they can be applied to any conforming approximation of the problem. However, from the practical viewpoint these estimates have an essential drawback:

they use a function τ in the set \mathbf{Q}_f defined by the differential relation,

which may be difficult to satisfy exactly. Probably by this reason further development of a posteriori error estimates for Finite Element Methods (especially in 80'-90') was mainly based on different grounds.

Fixed point theorem

Consider a Banach space (\mathbf{X}, \mathbf{d}) and a continuous operator

$$\mathfrak{T} : \mathbf{X} \rightarrow \mathbf{X}.$$

Definition

A point \mathbf{x}_\odot is called a fixed point of \mathfrak{T} if

$$\mathbf{x}_\odot = \mathfrak{T}\mathbf{x}_\odot. \quad (10)$$

Approximations of a fixed point are usually constructed by the iteration sequence

$$\mathbf{x}_i = \mathfrak{T}\mathbf{x}_{i-1} \quad i = 1, 2, \dots \quad (11)$$

Contractive mappings

Two basic tasks:

- (a) find the conditions that guarantee convergence of \mathbf{x}_i to \mathbf{x}_\odot ,
- (b) find computable estimates of the error $\mathbf{e}_i = \mathbf{d}(\mathbf{x}_i, \mathbf{x}_\odot)$.

Definition

An operator $\mathfrak{T} : \mathbf{X} \rightarrow \mathbf{X}$ is called q -contractive on a set $\mathbf{S} \subset \mathbf{X}$ if there exists a positive real number q such that the inequality

$$\mathbf{d}(\mathfrak{T}\mathbf{x}, \mathfrak{T}\mathbf{y}) \leq q \mathbf{d}(\mathbf{x}, \mathbf{y}) \quad (12)$$

holds for any elements \mathbf{x} and \mathbf{y} of the set \mathbf{S} .

Theorem (S. Banach)

Let \mathfrak{T} be a \mathbf{q} -contractive mapping of a closed nonempty set $\mathbf{S} \subset \mathbf{X}$ to itself with $\mathbf{q} < \mathbf{1}$. Then, \mathfrak{T} has a unique fixed point in \mathbf{S} and the sequence \mathbf{x}_i obtained by (11) converges to this point.

Proof. It is easy to see that

$$\mathbf{d}(\mathbf{x}_{i+1}, \mathbf{x}_i) = \mathbf{d}(\mathfrak{T}\mathbf{x}_i, \mathfrak{T}\mathbf{x}_{i-1}) \leq \mathbf{q}\mathbf{d}(\mathbf{x}_i, \mathbf{x}_{i-1}) \leq \dots \leq \mathbf{q}^i \mathbf{d}(\mathbf{x}_1, \mathbf{x}_0).$$

Therefore, for any $\mathbf{m} > \mathbf{1}$ we have

$$\begin{aligned} \mathbf{d}(\mathbf{x}_{i+m}, \mathbf{x}_i) &\leq \\ &\leq \mathbf{d}(\mathbf{x}_{i+m}, \mathbf{x}_{i+m-1}) + \mathbf{d}(\mathbf{x}_{i+m-1}, \mathbf{x}_{i+m-2}) + \dots + \mathbf{d}(\mathbf{x}_{i+1}, \mathbf{x}_i) \leq \\ &\leq \mathbf{q}^i (\mathbf{q}^{m-1} + \mathbf{q}^{m-2} + \dots + \mathbf{1}) \mathbf{d}(\mathbf{x}_1, \mathbf{x}_0). \quad (13) \end{aligned}$$

Since

$$\sum_{k=0}^{m-1} \mathbf{q}^k \leq \frac{\mathbf{1}}{\mathbf{1} - \mathbf{q}},$$

(13) implies the estimate

$$\mathbf{d}(\mathbf{x}_{i+m}, \mathbf{x}_i) \leq \frac{\mathbf{q}^i}{\mathbf{1} - \mathbf{q}} \mathbf{d}(\mathbf{x}_1, \mathbf{x}_0). \quad (14)$$

Let $\mathbf{i} \rightarrow \infty$, then the right-hand side of (14) tends to zero, so that $\{\mathbf{x}_i\}$ is a Cauchy sequence. It has a limit in $\mathbf{y} \in \mathbf{X}$.

Then, $\mathbf{d}(\mathbf{x}_i, \mathbf{y}) \rightarrow \mathbf{0}$ and

$$\mathbf{d}(\mathfrak{T}\mathbf{x}_i, \mathfrak{T}\mathbf{y}) \leq \mathbf{q}\mathbf{d}(\mathbf{x}_i, \mathbf{y}) \rightarrow \mathbf{0}$$

so that $\mathbf{d}(\mathfrak{T}\mathbf{x}_i, \mathfrak{T}\mathbf{y}) \rightarrow \mathbf{0}$ and $\mathfrak{T}\mathbf{x}_i \rightarrow \mathfrak{T}\mathbf{y}$. Pass to the limit in (11) as $i \rightarrow +\infty$. We observe that

$$\mathfrak{T}\mathbf{y} = \mathbf{y}.$$

Hence, **any limit of such a sequence is a fixed point.**

It is easy to prove that a fixed point is **unique**.

Assume that there are two different fixed points \mathbf{x}_{\odot}^1 and \mathbf{x}_{\odot}^2 , i.e.

$$\mathfrak{T}\mathbf{x}_{\odot}^k = \mathbf{x}_{\odot}^k, \quad \mathbf{k} = 1, 2.$$

Therefore,

$$\mathbf{d}(\mathbf{x}_{\odot}^1, \mathbf{x}_{\odot}^2) = \mathbf{d}(\mathfrak{T}\mathbf{x}_{\odot}^1, \mathfrak{T}\mathbf{x}_{\odot}^2) \leq \mathbf{q}\mathbf{d}(\mathbf{x}_{\odot}^1, \mathbf{x}_{\odot}^2).$$

But $\mathbf{q} < \mathbf{1}$, and thus such an inequality cannot be true.

A priori convergence estimate

Let $\mathbf{e}_j = \mathbf{d}(\mathbf{x}_j, \mathbf{x}_\odot)$ denote the error on the j -th step. Then

$$\mathbf{e}_j = \mathbf{d}(\mathfrak{T}\mathbf{x}_{j-1}, \mathfrak{T}\mathbf{x}_\odot) \leq \mathbf{q}\mathbf{e}_{j-1} \leq \mathbf{q}^j\mathbf{e}_0.$$

and

$$\mathbf{e}_j \leq \mathbf{q}^j\mathbf{e}_0. \quad (15)$$

This estimate gives a certain presentation on that how the error decreases. However, this a priori upper bound may be rather coarse.

A posteriori estimates

The proposition below furnishes upper and lower estimates of \mathbf{e}_j , which are easy to compute provided, that the number \mathbf{q} (or a good estimate of it) is known.

A. Ostrowski. Les estimations des erreurs a posteriori dans les procédés itératifs, *C.R. Acad.Sci. Paris Sér. A-B*, 275(1972), A275-A278.

Theorem (A. Ostrowski)

Let $\{\mathbf{x}_j\}_{j=0}^{\infty}$ be a sequence obtained by the iteration process

$$\mathbf{x}_i = \mathfrak{T}\mathbf{x}_{i-1} \quad i = 1, 2, \dots$$

with a mapping \mathfrak{T} satisfying the condition $\|\mathfrak{T}\| = \mathbf{q} \leq \mathbf{1}$. Then, for any \mathbf{x}_j , $j > \mathbf{1}$, the following estimate holds:

$$\mathbf{M}_{\ominus}^j := \frac{\mathbf{1}}{\mathbf{1} + \mathbf{q}} \mathbf{d}(\mathbf{x}_{j+1}, \mathbf{x}_j) \leq \mathbf{e}_j \leq \mathbf{M}_{\oplus}^j := \frac{\mathbf{q}}{\mathbf{1} - \mathbf{q}} \mathbf{d}(\mathbf{x}_j, \mathbf{x}_{j-1}). \quad (16)$$

Proof. The upper estimate in (16) follows from (14). Indeed, put $\mathbf{i} = \mathbf{1}$ in this relation. We have

$$\mathbf{d}(\mathbf{x}_{1+m}, \mathbf{x}_1) \leq \frac{\mathbf{q}}{\mathbf{1} - \mathbf{q}} \mathbf{d}(\mathbf{x}_1, \mathbf{x}_0).$$

Since $\mathbf{x}_{1+m} \rightarrow \mathbf{x}_\odot$ as $\mathbf{m} \rightarrow +\infty$, we pass to the limit with respect to \mathbf{m} and obtain

$$\mathbf{d}(\mathbf{x}_\odot, \mathbf{x}_1) \leq \frac{\mathbf{q}}{\mathbf{1} - \mathbf{q}} \mathbf{d}(\mathbf{x}_1, \mathbf{x}_0).$$

We may view \mathbf{x}_{j-1} as the starting point of the sequence. Then, in the above relation $\mathbf{x}_0 = \mathbf{x}_{j-1}$ and $\mathbf{x}_1 = \mathbf{x}_j$ and we arrive at the following **upper bound** of the error:

$$\mathbf{d}(\mathbf{x}_\odot, \mathbf{x}_j) \leq \frac{\mathbf{q}}{\mathbf{1} - \mathbf{q}} \mathbf{d}(\mathbf{x}_j, \mathbf{x}_{j-1}).$$

The **lower bound** of the error follows from the relation

$$\mathbf{d}(\mathbf{x}_j, \mathbf{x}_{j-1}) \leq \mathbf{d}(\mathbf{x}_j, \mathbf{x}_\odot) + \mathbf{d}(\mathbf{x}_{j-1}, \mathbf{x}_\odot) \leq (\mathbf{1} + \mathbf{q})\mathbf{d}(\mathbf{x}_{j-1}, \mathbf{x}_\odot),$$

which shows that

$$\mathbf{d}(\mathbf{x}_{j-1}, \mathbf{x}_\odot) \geq \frac{\mathbf{1}}{\mathbf{1} + \mathbf{q}} \mathbf{d}(\mathbf{x}_j, \mathbf{x}_{j-1}).$$

Note that

$$\frac{\mathbf{M}_\oplus^j}{\mathbf{M}_\ominus^j} = \frac{\mathbf{q}(\mathbf{1} + \mathbf{q})}{\mathbf{1} - \mathbf{q}} \frac{\mathbf{d}(\mathbf{x}_j, \mathbf{x}_{j-1})}{\mathbf{d}(\mathbf{x}_{j+1}, \mathbf{x}_j)} \geq \frac{\mathbf{1} + \mathbf{q}}{\mathbf{1} - \mathbf{q}},$$

we see that that the efficiency of the upper and lower bounds given by (16) deteriorates as $\mathbf{q} \rightarrow \mathbf{1}$.

Remark. If \mathbf{X} is a normed space, then

$$\mathbf{d}(\mathbf{x}_{j+1}, \mathbf{x}_j) = \|\mathbf{R}(\mathbf{x}_j)\|,$$

where

$$\mathbf{R}(\mathbf{x}_j) := \mathfrak{T}\mathbf{x}_j - \mathbf{x}_j$$

is the residual of the basic equation (10). Thus, the upper and lower estimates of errors are expressed in terms of the **residuals of the respective iteration equation** computed for two neighbor steps:

$$\frac{1}{1+q} \|\mathbf{R}(\mathbf{x}_j)\| \leq \mathbf{e}_j = \mathbf{d}(\mathbf{x}_j, \mathbf{x}_\odot) \leq \frac{q}{1-q} \|\mathbf{R}(\mathbf{x}_{j-1})\|.$$

Corollaries

In the iteration methods, it is often easier to analyze the operator

$$\mathfrak{T} = \mathbf{T}^n := \underbrace{\mathbf{T}\mathbf{T}\dots\mathbf{T}}_{n \text{ times}}$$

where \mathbf{T} is a certain mapping.

Proposition (1)

Let $\mathbf{T} : \mathbf{S} \rightarrow \mathbf{S}$ be a continuous mapping such that \mathfrak{T} is a \mathbf{q} -contractive mapping with $\mathbf{q} \in (0, 1)$. Then, the equations

$$\mathbf{x} = \mathbf{T}\mathbf{x} \quad \text{and} \quad \mathbf{x} = \mathfrak{T}\mathbf{x}$$

have one and the same fixed point, which is unique and can be found by the above described iteration procedure.

Proof. By the Banach Theorem, we observe that the operator \mathfrak{T} has a unique fixed point ξ_{\odot} .

Let us show that ξ_{\odot} is a fixed point of \mathbf{T} , we note that

$$\begin{aligned} \mathbf{T}\xi_{\odot} &= \mathbf{T}(\mathfrak{T}\xi_{\odot}) = \mathbf{T}\mathfrak{T}^2\xi_{\odot} = \dots \\ &= \mathbf{T}\mathfrak{T}^i\xi_{\odot} = \mathbf{T}^{(1+i\text{in})}\xi_{\odot} = \mathbf{T}^{\text{in}}\mathbf{T}\xi_{\odot}. \end{aligned} \quad (17)$$

Denote $\mathbf{x}_0 = \mathbf{T}\xi_{\odot}$. By (17) we conclude that for any i

$$\mathbf{T}\xi_{\odot} = \mathfrak{T}^i\mathbf{x}_0. \quad (18)$$

Passing to the limit on the right-hand side in (18), we arrive at the relation $\mathbf{T}\xi_{\odot} = \xi_{\odot}$, which means that ξ_{\odot} is a fixed point of the operator \mathbf{T} .

Let $\widetilde{\mathbf{x}}_{\odot}$ be a fixed point of \mathbf{T} . Then,

$$\widetilde{\mathbf{x}}_{\odot} = \mathbf{T}^2 \widetilde{\mathbf{x}}_{\odot} = \dots = \mathbf{T}^n \widetilde{\mathbf{x}}_{\odot} = \mathfrak{T} \widetilde{\mathbf{x}}_{\odot}$$

and we observe that $\widetilde{\mathbf{x}}_{\odot}$ is a fixed point of \mathbf{T} . Since the saddle point of \mathfrak{T} exists and is unique, we conclude that

$$\mathbf{x}_{\odot} = \widetilde{\mathbf{x}}_{\odot}.$$

Remark. This assertion may be practically useful if it is not possible to prove that \mathbf{T} is \mathbf{q} -contractive, but this fact can be established for a certain power of \mathbf{T} .

Iteration methods for bounded linear operators

Consider a bounded linear operator $\mathcal{L} : \mathbf{X} \rightarrow \mathbf{X}$, where \mathbf{X} is a Banach space. Given $\mathbf{b} \in \mathbf{X}$, the iteration process is defined by the relation

$$\mathbf{x}_j = \mathcal{L} \mathbf{x}_{j-1} + \mathbf{b}. \quad (19)$$

Let \mathbf{x}_\odot be a fixed point of (19) and

$$\|\mathcal{L}\| = \mathbf{q} < \mathbf{1}.$$

By applying the Banach Theorem it is easy to show that

$$\{\mathbf{x}_j\} \rightarrow \mathbf{x}_\odot.$$

Indeed, let $\bar{\mathbf{x}}_j = \mathbf{x}_j - \mathbf{x}_\odot$. Then

$$\bar{\mathbf{x}}_j = \mathcal{L}\mathbf{x}_{j-1} + \mathbf{b} - \mathbf{x}_\odot = \mathcal{L}(\mathbf{x}_{j-1} - \mathbf{x}_\odot) = \mathcal{L}\bar{\mathbf{x}}_{j-1}. \quad (20)$$

Since

$$\mathbf{0}_X = \mathcal{L}\mathbf{0}_X,$$

we note that the zero element $\mathbf{0}_X$ is a unique fixed point of the operator \mathcal{L} . By the Banach theorem $\bar{\mathbf{x}}_j \rightarrow \mathbf{0}_X$ and, therefore, $\{\mathbf{x}_j\} \rightarrow \mathbf{x}_\odot$.

Therefore, we have an *a priori* estimate

$$\begin{aligned}\|\mathbf{x}_j - \mathbf{x}_\odot\|_{\mathbf{X}} &= \|\bar{\mathbf{x}}_j - \mathbf{0}_{\mathbf{X}}\|_{\mathbf{X}} \leq \\ &\leq \frac{\mathbf{q}^j}{\mathbf{1} - \mathbf{q}} \|\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_0\|_{\mathbf{X}} = \frac{\mathbf{q}^j}{\mathbf{1} - \mathbf{q}} \|\mathbf{R}(\mathbf{x}_0)\|_{\mathbf{X}}\end{aligned}\quad (21)$$

and the *a posteriori* one

$$\|\mathbf{x}_j - \mathbf{x}_\odot\|_{\mathbf{X}} \leq \frac{\mathbf{q}}{\mathbf{1} - \mathbf{q}} \|\mathbf{R}(\mathbf{x}_{j-1})\|_{\mathbf{X}},\quad (22)$$

where $\mathbf{R}(\mathbf{z}) = \mathcal{L}\mathbf{z} + \mathbf{b} - \mathbf{z}$ is the *residual* of the functional equation considered.

By applying the general theory, we also obtain a lower bound of the error

$$\|\mathbf{x}_j - \mathbf{x}_\odot\|_{\mathbf{X}} \geq \frac{\mathbf{1}}{\mathbf{1} + \mathbf{q}} \|\mathbf{x}_{j+1} - \mathbf{x}_j\|_{\mathbf{X}} = \frac{\mathbf{1}}{\mathbf{1} + \mathbf{q}} \|\mathbf{R}(\mathbf{x}_j)\|_{\mathbf{X}}. \quad (23)$$

Hence, we arrive at the following estimates for the error in the linear operator equation:

$$\frac{\mathbf{1} - \mathbf{q}}{\mathbf{q}} \|\mathbf{x}_j - \mathbf{x}_\odot\|_{\mathbf{X}} \leq \|\mathbf{R}(\mathbf{x}_{j-1})\|_{\mathbf{X}} \leq (\mathbf{1} + \mathbf{q}) \|\mathbf{x}_{j-1} - \mathbf{x}_\odot\|_{\mathbf{X}}.$$

Iteration methods in linear algebra

Important applications of the above results are associated with systems of linear simultaneous equations and other algebraic problems. Set $\mathbf{X} = \mathbb{R}^n$ and assume that \mathcal{L} is defined by a nondegenerate matrix $\mathbf{A} \in \mathbb{M}^{n \times n}$ decomposed into three matrixes

$$\mathbf{A} = \mathbf{A}_\ell + \mathbf{A}_d + \mathbf{A}_r,$$

where \mathbf{A}_ℓ , \mathbf{A}_r , and \mathbf{A}_d are certain lower, upper, and diagonal matrices, respectively.

Iteration methods for systems of linear simultaneous equations associated with \mathbf{A} are often represented in the form

$$\mathbf{B} \frac{\mathbf{x}_i - \mathbf{x}_{i-1}}{\tau} + \mathbf{A} \mathbf{x}_{i-1} = \mathbf{f}. \quad (24)$$

In (24), the matrix \mathbf{B} and the parameter τ may be taken in various ways (depending on the properties of \mathbf{A}). We consider three frequently encountered cases:

- (a) $\mathbf{B} = \mathbf{A}_d$,
- (b) $\mathbf{B} = \mathbf{A}_d + \mathbf{A}_\ell$,
- (c) $\mathbf{B} = \mathbf{A}_d + \omega \mathbf{A}_\ell$, $\tau = \omega$.

For $\tau = 1$, (a) and (b) lead to the methods of Jacobi and Zeidel, respectively. In (c), the parameter ω must be in the interval $(0, 2)$. If $\omega > 1$, we have the so-called "upper relaxation method", and $\omega < 1$ corresponds to the "lower relaxation method".

The method (24) is reduced to (19) if we set

$$\mathcal{L} = \mathbb{I} - \tau \mathbf{B}^{-1} \mathbf{A} \quad \text{and} \quad \mathbf{b} = \tau \mathbf{B}^{-1} \mathbf{f}, \quad (25)$$

where \mathbb{I} is the unit matrix. It is known that \mathbf{x}_i converges to \mathbf{x}_\odot that is a solution of the system

$$\mathbf{A} \mathbf{x}_\odot = \mathbf{f} \quad (26)$$

if and only if all the eigenvalues of \mathcal{L} are less than one.

Obviously, \mathbf{B} and τ should be taken in such a way that they guarantee the fulfillment of this condition.

Assume that $\|\mathcal{L}\| \leq \mathbf{q} < \mathbf{1}$. In view of (21)-(23), the quantities

$$\mathbf{M}_{\oplus}^i = \mathbf{q}(\mathbf{1} - \mathbf{q})^{-1} \|\mathbf{R}(\mathbf{x}_{i-1})\|, \quad (27)$$

$$\mathbf{M}_{\oplus}^{0i} = \mathbf{q}^i(\mathbf{1} - \mathbf{q})^{-1} \|\mathbf{R}(\mathbf{x}_0)\|, \quad (28)$$

$$\mathbf{M}_{\ominus}^i = (\mathbf{1} + \mathbf{q})^{-1} \|\mathbf{R}(\mathbf{x}_i)\| \quad (29)$$

furnish upper and lower bounds of the error for the vector \mathbf{x}_i .

Remark. It is worth noting that from the practical viewpoint finding an upper bound for $\|\mathcal{L}\|$ and proving that it is less than 1 presents a special and often not easy task.

If \mathbf{q} is very close to 1, then the convergence of an iteration process may be very slow. As we have seen, in this case, the quality of error estimates is also degraded. A well-accepted way for accelerating the convergence consists of using a modified system obtained from the original one by means of a suitable *preconditioner*

Task 2

Consider the problem

$$\mathbf{Ax} = \mathbf{f}$$

for a symmetric matrix \mathbf{A} with coefficients

$$\mathbf{a}_{ij} = \kappa/ij \quad \text{if } i \neq j, \quad \kappa = 0.1$$

$$\mathbf{a}_{ii} = i.$$

Solved the system by the iteration method

$$\mathbf{x}_{i+1} = (\mathbb{I} - \tau \mathbf{B}^{-1} \mathbf{A}) \mathbf{x}_i + \tau \mathbf{B}^{-1} \mathbf{F}$$

with $\mathbf{B} = \mathbf{A}_D$ and $\mathbf{x}_0 = \{\mathbf{0}, \mathbf{0}, \dots, \mathbf{0}\}$, determine q and define two-sided error bounds.

Applications to integral equations

Many problems in science and engineering can be stated in terms of integral equations. One of the most typical cases is to find a function $\mathbf{x}_\odot(\mathbf{t}) \in \mathbf{C}[\mathbf{a}, \mathbf{b}]$ such that

$$\mathbf{x}_\odot(\mathbf{t}) = \lambda \int_{\mathbf{a}}^{\mathbf{b}} \mathbf{K}(\mathbf{t}, \mathbf{s}) \mathbf{x}_\odot(\mathbf{s}) \, \mathbf{d}\mathbf{s} + \mathbf{f}(\mathbf{t}), \quad (30)$$

where $\lambda \geq 0$, \mathbf{K} (the kernel) is a continuous function for

$$(\mathbf{x}, \mathbf{t}) \in \mathbf{Q} := \{\mathbf{a} \leq \mathbf{s} \leq \mathbf{b}, \mathbf{a} \leq \mathbf{t} \leq \mathbf{b}\}$$

and

$$|\mathbf{K}(\mathbf{t}, \mathbf{s})| \leq \mathbf{M}, \quad \forall (\mathbf{t}, \mathbf{s}) \in \mathbf{Q}.$$

Also, we assume that $\mathbf{f} \in \mathbf{C}[\mathbf{a}, \mathbf{b}]$.

Let us define the operator \mathfrak{T} as follows:

$$\mathbf{y}(\mathbf{t}) := \mathfrak{T}\mathbf{x}(\mathbf{t}) := \lambda \int_a^b \mathbf{K}(\mathbf{t}, \mathbf{x})\mathbf{x}(\mathbf{s}) \, d\mathbf{s} + \mathbf{f}(\mathbf{t}) \quad (31)$$

and show that \mathfrak{T} maps continuous functions to continuous ones. Let \mathbf{t}_0 and $\mathbf{t}_0 + \Delta\mathbf{t}$ belong to $[\mathbf{a}, \mathbf{b}]$. Then,

$$\begin{aligned} |\mathbf{y}(\mathbf{t}_0 + \Delta\mathbf{t}) - \mathbf{y}(\mathbf{t}_0)| &\leq \\ &\leq |\lambda| \int_a^b |\mathbf{K}(\mathbf{t}_0 + \Delta\mathbf{t}, \mathbf{s}) - \mathbf{K}(\mathbf{t}_0, \mathbf{s})| |\mathbf{x}(\mathbf{s})| \, d\mathbf{s} + \\ &\quad + |\mathbf{f}(\mathbf{t}_0 + \Delta\mathbf{t}) - \mathbf{f}(\mathbf{t}_0)|. \end{aligned}$$

Since \mathbf{K} and \mathbf{f} are continuous on the compact sets \mathbf{Q} and $[\mathbf{a}, \mathbf{b}]$, respectively, they are uniformly continuous on these sets.

Therefore, for any given ε one can find a small number δ such that

$$|\mathbf{f}(\mathbf{t}_0 + \Delta \mathbf{t}) - \mathbf{f}(\mathbf{t}_0)| < \varepsilon$$

and

$$|\mathbf{K}(\mathbf{t}_0 + \Delta \mathbf{t}, \mathbf{s}) - \mathbf{K}(\mathbf{t}_0, \mathbf{s})| < \varepsilon,$$

provided that $|\Delta \mathbf{t}| < \delta$.

Thus, we have

$$|\mathbf{y}(\mathbf{t}_0 + \Delta \mathbf{t}) - \mathbf{y}(\mathbf{t}_0)| \leq \varepsilon(|\lambda| \|\mathbf{b} - \mathbf{a}\| \max_{\mathbf{s} \in [\mathbf{a}, \mathbf{b}]} |\mathbf{x}(\mathbf{s})| + 1) = \mathbf{C}\varepsilon,$$

and, consequently, $\mathbf{y}(\mathbf{t}_0 + \Delta \mathbf{t})$ tends to $\mathbf{y}(\mathbf{t}_0)$ as $|\Delta \mathbf{t}| \rightarrow 0$.

$\mathfrak{T} : \mathbf{C}[a, b] \rightarrow \mathbf{C}[a, b]$ is a **contractive mapping**. Indeed,

$$\begin{aligned} d(\mathfrak{T}\mathbf{x}, \mathfrak{T}\mathbf{y}) &= \max_{a \leq t \leq b} |\mathfrak{T}\mathbf{x}(t) - \mathfrak{T}\mathbf{y}(t)| = \\ &= \max_{a \leq t \leq b} \left| \lambda \int_a^b \mathbf{K}(t, s)(\mathbf{x}(s) - \mathbf{y}(s)) \, ds \right| \leq \\ &\leq |\lambda| \mathbf{M}(\mathbf{b} - \mathbf{a}) \max_{a \leq s \leq b} |\mathbf{x}(s) - \mathbf{y}(s)| = |\lambda| \mathbf{M}(\mathbf{b} - \mathbf{a}) d(\mathbf{x}, \mathbf{y}), \end{aligned}$$

so that \mathfrak{T} is a **q**-contractive operator with

$$\mathbf{q} = |\lambda| \mathbf{M}(\mathbf{b} - \mathbf{a}), \quad (32)$$

provided that

$$|\lambda| < \frac{1}{\mathbf{M}(\mathbf{b} - \mathbf{a})}. \quad (33)$$

Numerical procedure

An approximate solution of (30) can be found by the iteration method

$$\mathbf{x}_{i+1}(\mathbf{t}) = \lambda \int_a^b \mathbf{K}(\mathbf{t}, \mathbf{s}) \mathbf{x}_i(\mathbf{s}) \, d\mathbf{s} + \mathbf{f}(\mathbf{t}). \quad (34)$$

If (33) holds, then from the Banach theorem it follows that the sequence $\{\mathbf{x}_i\}$ converges to the exact solution.

We apply the theory exposed above and find that the accuracy of \mathbf{x}_i is subject to the estimate

$$\begin{aligned} \frac{1}{1+q} \int_a^b \mathbf{K}(\mathbf{t}, \mathbf{s}) (\mathbf{x}_{i+1}(\mathbf{s}) - \mathbf{x}_i(\mathbf{s})) \, d\mathbf{s} &\leq \\ &\leq \max_{a \leq \mathbf{t} \leq b} |\mathbf{x}_i(\mathbf{t}) - \mathbf{x}_\odot(\mathbf{t})| \leq \frac{q}{1-q} \int_a^b \mathbf{K}(\mathbf{t}, \mathbf{s}) (\mathbf{x}_i(\mathbf{s}) - \mathbf{x}_{i-1}(\mathbf{s})) \, d\mathbf{s}. \end{aligned} \quad (35)$$

Applications to Volterra type equations

Consider the fixed point problem

$$\mathbf{x}_{\odot}(\mathbf{t}) = \lambda \int_a^{\mathbf{t}} \mathbf{K}(\mathbf{t}, \mathbf{s}) \mathbf{x}_{\odot}(\mathbf{s}) \, \mathbf{d}\mathbf{s} + \mathbf{f}(\mathbf{t}), \quad (36)$$

where

$$|\mathbf{K}(\mathbf{t}, \mathbf{s})| \leq \mathbf{M}, \quad \forall (\mathbf{t}, \mathbf{s}) \in \mathbf{Q}$$

and $\mathbf{f} \in \mathbf{C}[\mathbf{a}, \mathbf{b}]$.

Define the operator \mathbf{T} as follows:

$$\mathbf{T}\mathbf{x}(\mathbf{t}) = \lambda \int_a^{\mathbf{t}} \mathbf{K}(\mathbf{t}, \mathbf{s}) \mathbf{x}(\mathbf{s}) \, \mathbf{d}\mathbf{s} + \mathbf{f}(\mathbf{t}).$$

Similarly, to the previous case we establish that

$$\mathbf{d}(\mathbf{T}\mathbf{x}, \mathbf{T}\mathbf{y}) \leq |\lambda| \mathbf{M}(\mathbf{t} - \mathbf{a}) \mathbf{d}(\mathbf{x}, \mathbf{y}).$$

By the same arguments we find that

$$\mathbf{d}(\mathbf{T}^n \mathbf{x}, \mathbf{T}^n \mathbf{y}) \leq |\lambda|^n \mathbf{M}^n \frac{(\mathbf{t} - \mathbf{a})^n}{\mathbf{n}!} \mathbf{d}(\mathbf{x}, \mathbf{y}),$$

Thus, the operator $\mathfrak{T} := \mathbf{T}^n$ is \mathbf{q} -contractive with a certain $\mathbf{q} < \mathbf{1}$, provided that \mathbf{n} is large enough.

In view of Proposition 1, we conclude that the iteration method converges to \mathbf{x}_\odot and the errors are controlled by the two-sided error estimates.

Applications to ordinary differential equations

Let \mathbf{u} be a solution of the simplest initial boundary-value problem

$$\frac{d\mathbf{u}}{dt} = \varphi(\mathbf{t}, \mathbf{u}(\mathbf{t})), \quad \mathbf{u}(\mathbf{t}_0) = \mathbf{a}, \quad (37)$$

where the solution $\mathbf{u}(\mathbf{t})$ is to be found on the interval $[\mathbf{t}_0, \mathbf{t}_1]$. Assume that the function $\varphi(\mathbf{t}, \mathbf{p})$ is continuous on the set

$$\mathbf{Q} = \{\mathbf{t}_0 \leq \mathbf{t} \leq \mathbf{t}_1, \mathbf{a} - \Delta \leq \mathbf{p} \leq \mathbf{a} + \Delta\}$$

and

$$|\varphi(\mathbf{t}, \mathbf{p}_1) - \varphi(\mathbf{t}, \mathbf{p}_2)| \leq \mathbf{L}|\mathbf{p}_1 - \mathbf{p}_2|, \quad \forall(\mathbf{t}, \mathbf{p}) \in \mathbf{Q}. \quad (38)$$

Problem (37) can be reduced to the integral equation

$$\mathbf{u}(\mathbf{t}) = \int_{t_0}^{\mathbf{t}} \varphi(\mathbf{s}, \mathbf{u}(\mathbf{s})) \, d\mathbf{s} + \mathbf{a} \quad (39)$$

and it is natural to solve the latter problem by the iteration method

$$\mathbf{u}_j(\mathbf{t}) = \int_{t_0}^{\mathbf{t}} \varphi(\mathbf{s}, \mathbf{u}_{j-1}(\mathbf{s})) \, d\mathbf{s} + \mathbf{a}. \quad (40)$$

To justify this procedure, we must verify that the operator

$$\mathfrak{T}\mathbf{u} := \int_{t_0}^{\mathbf{t}} \varphi(\mathbf{s}, \mathbf{u}(\mathbf{s})) \, d\mathbf{s} + \mathbf{a}$$

is \mathbf{q} -contractive with respect to the norm

$$\|\mathbf{u}\| := \max_{\mathbf{t} \in [t_0, t_1]} |\mathbf{u}(\mathbf{t})|. \quad (41)$$

We have

$$\begin{aligned} \|\mathfrak{T}z - \mathfrak{T}y\| &= \max_{t \in [t_0, t_1]} \left| \int_{t_0}^t (\varphi(s, z(s)) - \varphi(s, y(s))) ds \right| \leq \\ &\leq \max_{t \in [t_0, t_1]} L \int_{t_0}^t |z(s) - y(s)| ds \leq L \int_{t_0}^{t_1} |z(s) - y(s)| ds \leq \\ &\leq L(t_1 - t_0) \max_{s \in [t_0, t_1]} |z(s) - y(s)| = L(t_1 - t_0) \|z - y\|. \end{aligned}$$

We see that if

$$t_1 < t_0 + L^{-1}, \quad (42)$$

then the operator \mathfrak{T} is \mathbf{q} -contractive with

$$\mathbf{q} := L(t_1 - t_0) < 1.$$

Therefore, if the interval $[t_0, t_1]$ is small enough (i.e., it satisfies the condition 42), then the existence and uniqueness of a continuous solution $\mathbf{u}(\mathbf{t})$ follows from the Banach theorem. In this case, the solution can be found by the iteration procedure whose accuracy is explicitly controlled by the two-sided error estimates.

For a more detailed investigation of the fixed point methods for integral and differential equations see

A. N. Kolmogorov and S. V. Fomin. *Introductory real analysis*. Dover Publications, Inc., New York, 1975.

E. Zeidler. *Nonlinear functional analysis and its applications. I. Fixed-point theorems*. Springer-Verlag, New York, 1986.

A posteriori estimates based on monotonicity.

The theory of **monotone operators** gives another way of constructing a posteriori estimates.

Monotone operators are defined on the so-called **ordered** (or **partially ordered**) spaces that introduce the relation $\mathbf{x} \leq \mathbf{y}$ for all (or almost all) elements \mathbf{x}, \mathbf{y} of the space.

Definition

An operator \mathfrak{T} is called monotone if $\mathbf{x} \leq \mathbf{y}$ implies $\mathfrak{T}\mathbf{x} \leq \mathfrak{T}\mathbf{y}$.

Consider the fixed point problem

$$\mathbf{x}_\ominus = \mathfrak{T}\mathbf{x}_\ominus + \mathbf{f}$$

on an ordered (partially ordered) space X . Assume that

$$\mathfrak{T} = \mathfrak{T}_\oplus + \mathfrak{T}_\ominus,$$

\mathfrak{T}_\oplus is monotone,

\mathfrak{T}_\ominus is antitone: $\mathbf{x} \leq \mathbf{y}$ implies $\mathfrak{T}\mathbf{x} \geq \mathfrak{T}\mathbf{y}$,

\mathfrak{T}_\oplus and \mathfrak{T}_\ominus have a common set of images \mathbf{D} which is a convex subset of \mathbf{X} .

Next, let $\mathbf{x}_{\ominus 0}, \mathbf{x}_{\ominus 1}, \mathbf{x}_{\oplus 0}, \mathbf{x}_{\oplus 1} \in \mathbf{D}$ be such elements that

$$\begin{aligned}\mathbf{x}_{\ominus 0} &\leq \mathbf{x}_{\ominus 1} \leq \mathbf{x}_{\oplus 1} \leq \mathbf{x}_{\oplus 0}, \\ \mathbf{x}_{\ominus 1} &= \mathfrak{T}_{\oplus} \mathbf{x}_{\ominus 0} + \mathfrak{T}_{\ominus} \mathbf{x}_{\oplus 0} + \mathbf{f}, \\ \mathbf{x}_{\oplus 1} &= \mathfrak{T}_{\oplus} \mathbf{x}_{\oplus 0} + \mathfrak{T}_{\ominus} \mathbf{x}_{\ominus 0} + \mathbf{f},\end{aligned}$$

Then, we observe that

$$\begin{aligned}\mathbf{x}_{\ominus 2} &= \mathfrak{T}_{\oplus} \mathbf{x}_{\ominus 1} + \mathfrak{T}_{\ominus} \mathbf{x}_{\oplus 1} + \mathbf{f} \geq \mathfrak{T}_{\oplus} \mathbf{x}_{\ominus 0} + \mathfrak{T}_{\ominus} \mathbf{x}_{\oplus 0} + \mathbf{f} = \mathbf{x}_{\ominus 1} \\ \mathbf{x}_{\oplus 2} &= \mathfrak{T}_{\oplus} \mathbf{x}_{\oplus 1} + \mathfrak{T}_{\ominus} \mathbf{x}_{\ominus 1} + \mathbf{f} \leq \mathfrak{T}_{\oplus} \mathbf{x}_{\oplus 0} + \mathfrak{T}_{\ominus} \mathbf{x}_{\ominus 0} + \mathbf{f} = \mathbf{x}_{\oplus 1}.\end{aligned}$$

By continuing the iterations we obtain elements such that

$$\mathbf{x}_{\ominus k} \leq \mathbf{x}_{\ominus(k+1)} \leq \mathbf{x}_{\oplus(k+1)} \leq \mathbf{x}_{\oplus k}.$$

Then $\mathbf{x} \rightarrow \mathfrak{T}\mathbf{x} + \mathbf{f}$ maps \mathbf{D} to itself. If \mathbf{D} is compact, then by the Schauder fixed point theorem $\mathbf{x}_{\odot} \in \mathbf{D}$ exists. Moreover, it is bounded from below and above by the sequences $\{\mathbf{x}_{\ominus k}\}$ and $\{\mathbf{x}_{\oplus k}\}$.

Applications of this method are mainly oriented towards systems of linear simultaneous equations and integral equations

(see [L. Collatz. Funktionanalysis und numerische mathematik, Springer-Verlag, Berlin, 1964](#)). For example, consider a system of linear simultaneous equations

$$\mathbf{x} = \mathbf{Ax} + \mathbf{f}$$

that is supposed to have a unique solution \mathbf{x}_{\odot} . Assume that

$$\begin{aligned} \mathbf{A} &= \mathbf{A}_{\oplus} - \mathbf{A}_{\ominus}, & \mathbf{A}_{\ominus} &= \{\mathbf{a}_{ij}^{\ominus}\} \in \mathbb{M}^{n \times n}, \\ \mathbf{A}_{\oplus} &= \{\mathbf{a}_{ij}^{\oplus}\} \in \mathbb{M}^{n \times n}, & \mathbf{a}_{ij}^{\ominus} &\geq \mathbf{0}, \quad \mathbf{a}_{ij}^{\oplus} \geq \mathbf{0}. \end{aligned}$$

We may [partially order](#) the space \mathbb{R}^n by saying that $\mathbf{x} \leq \mathbf{y}$ if and only if $\mathbf{x}_i \leq \mathbf{y}_i$ for $i = 1, 2, \dots, n$. Compute the vectors

$$\mathbf{x}_{\ominus(k+1)} = \mathbf{A}_{\oplus} \mathbf{x}_{\ominus k} + \mathbf{A}_{\ominus} \mathbf{x}_{\oplus k} + \mathbf{f}, \quad \mathbf{x}_{\oplus(k+1)} = \mathbf{A}_{\oplus} \mathbf{x}_{\oplus k} + \mathbf{A}_{\ominus} \mathbf{x}_{\ominus k} + \mathbf{f}.$$

If $\mathbf{x}_{\ominus 0} \leq \mathbf{x}_{\ominus 1} \leq \mathbf{x}_{\odot} \leq \mathbf{x}_{\oplus 1} \leq \mathbf{x}_{\oplus 0}$, then for all the components of \mathbf{x}_{\odot} we obtain two-sided estimates

$$\mathbf{x}_{\ominus k}^{(i)} \leq \mathbf{x}_{\ominus(k+1)}^{(i)} \leq \mathbf{x}_{\odot}^{(i)} \leq \mathbf{x}_{\oplus(k+1)}^{(i)} \leq \mathbf{x}_{\oplus k}^{(i)}, \quad \mathbf{i} = 1, 2, \dots, n.$$

Task 3.

Apply the above method for finding two-sided bounds of the Euclid error norm and componentwise errors for a system of linear simultaneous equations

$$\mathbf{Ax} = \mathbf{f}$$

where

$$\begin{aligned} \mathbf{a}_{ij} &= (-1)^{i+j} \kappa / ij \quad \text{if } i \neq j, \quad \kappa = \mathbf{0.1} \\ \mathbf{a}_{ii} &= \mathbf{i}. \end{aligned}$$

For the i th component of the solution determine the lower and upper bounds as follows:

$$\max_{j=0,1,\dots,k+1} (\mathbf{x}_j^\ominus)_i \leq (\mathbf{x}_\odot)_i \leq \min_{j=0,1,\dots,k+1} (\mathbf{x}_j^\oplus)_i.$$

It should be remarked that convergence of $\mathbf{x}_{\ominus \mathbf{k}}^{(i)}$ and $\mathbf{x}_{\oplus \mathbf{k}}^{(i)}$ to \mathbf{x}_{\odot} (and the convergence rate) requires a special investigation, which must use specific features of a particular problem.

In principle, a posteriori error estimates based on monotonicity can provide the most informative POINTWISE a posteriori error estimates. Regrettably, the respective theory has not been yet properly investigated.

Lecture 2

The goal of Lecture 2 to give an overview of a posteriori error estimation methods developed for Finite Element approximations in 70th–80th.

Lecture plan

- **Mathematical background;**
- **Residual type error estimates;**
 - **Basic idea;**
 - **Estimates in 1D case;**
 - **Estimates in 2D case;**
 - **Comments;**
- **Methods based on post-processing;**
- **Methods using adjoint problems;**

Sobolev spaces with negative indices

Definition

Linear functionals defined on the functions of the space $\mathring{C}^\infty(\Omega)$ are called **distributions**. They form the space $\mathcal{D}'(\Omega)$

Value of a **distribution** \mathbf{g} on a function φ is $\langle \mathbf{g}, \varphi \rangle$.

Distributions possess an important property:

they have derivatives of any order

Let $\mathbf{g} \in \mathcal{D}'(\Omega)$, then the quantity $-\langle \mathbf{g}, \frac{\partial \varphi}{\partial x_i} \rangle$ is another linear functional on $\mathcal{D}(\Omega)$. It is viewed as a generalized partial derivative of \mathbf{g} taken over the i -th variable.

Derivatives of L^q -functions

Any function g from the space $L^q(\Omega)$ ($q \geq 1$) defines a certain distribution as

$$\langle \mathbf{g}, \varphi \rangle = \int_{\Omega} \mathbf{g} \varphi \, dx$$

and, therefore, has generalized derivatives of any order. The sets of distributions, which are derivatives of q -integrable functions, are called **Sobolev spaces with negative indices**.

Definition

The space $W^{-\ell,q}(\Omega)$ is the space of distributions $\mathbf{g} \in \mathcal{D}'(\Omega)$ such that

$$\mathbf{g} = \sum_{|\alpha| \leq \ell} \mathbf{D}^\alpha \mathbf{g}_\alpha,$$

where $\mathbf{g}_\alpha \in L^q(\Omega)$.

Spaces $W^{-1,p}(\Omega)$

$W^{-1,p}(\Omega)$ contains distributions that can be viewed as generalized derivatives of L^q -functions. The functional

$$\left\langle \frac{\partial \mathbf{f}}{\partial \mathbf{x}_i}, \varphi \right\rangle := - \int_{\Omega} \mathbf{f} \frac{\partial \varphi}{\partial \mathbf{x}_i} \, d\mathbf{x} \quad \mathbf{f} \in \mathbf{L}^q(\Omega)$$

is linear and continuous not only for $\varphi \in \mathring{C}^\infty(\Omega)$ but, also, for $\varphi \in \mathring{W}^{1,p}(\Omega)$, where $1/p + 1/q = 1$ (density property). Hence, first generalized derivatives of \mathbf{f} lie in the space dual to $\mathring{W}^{1,p}(\Omega)$ denoted by $\mathbf{W}^{-1,p}(\Omega)$.

For $\mathring{W}^{1,2}(\Omega) = \mathring{H}^1(\Omega)$, the respective dual space is denoted by $\mathbf{H}^{-1}(\Omega)$.

Norms in "negative spaces"

For $\mathbf{g} \in H^{-1}(\Omega)$ we may introduce two equivalent "negative norms".

$$\|\mathbf{g}\|_{(-1),\Omega} := \sup_{\varphi \in \mathring{H}^1(\Omega)} \frac{|\langle \mathbf{g}, \varphi \rangle|}{\|\varphi\|_{1,2,\Omega}} < +\infty$$

$$\|\mathbf{g}\| := \sup_{\varphi \in \mathring{H}^1(\Omega)} \frac{|\langle \mathbf{g}, \varphi \rangle|}{\|\nabla \varphi\|_{\Omega}} < +\infty$$

From the definitions, it follows that

$$\langle \mathbf{g}, \varphi \rangle \leq \|\mathbf{g}\|_{(-1),\Omega} \|\varphi\|_{1,2,\Omega}$$

$$\langle \mathbf{g}, \varphi \rangle \leq \|\mathbf{g}\| \|\nabla \varphi\|_{\Omega}$$

Errors and Residuals. First glance

If an analyst is not sure in the quality of an approximate solution computed, then the very first idea that comes to his mind is to substitute the approximate solution into the equation and look at the **equation residual**.

We begin by recalling basic relations between residuals and errors that hold for systems of **linear simultaneous equations**. Let $\mathcal{A} \in \mathbb{M}^{n \times n}$, $\det \mathcal{A} \neq 0$, consider the system

$$\mathcal{A}\mathbf{u} + \mathbf{f} = \mathbf{0}.$$

For any \mathbf{v} we have the simplest **residual** type estimate

$$\mathcal{A}(\mathbf{v} - \mathbf{u}) = \mathcal{A}\mathbf{v} + \mathbf{f}; \quad \Rightarrow \quad \|\mathbf{e}\| \leq \|\mathcal{A}^{-1}\| \|\mathbf{r}\|.$$

where $\mathbf{e} = \mathbf{v} - \mathbf{u}$ and $\mathbf{r} = \mathcal{A}\mathbf{v} + \mathbf{f}$.

Two-sided estimates

Define the quantities

$$\lambda_{\min} = \min_{\substack{y \in \mathbb{R}^n \\ y \neq 0}} \frac{\|\mathcal{A}y\|}{\|y\|} \quad \text{and} \quad \lambda_{\max} = \max_{\substack{y \in \mathbb{R}^n \\ y \neq 0}} \frac{\|\mathcal{A}y\|}{\|y\|}$$

Since $\mathcal{A}\mathbf{e} = \mathbf{r}$, we see that

$$\lambda_{\min} \leq \frac{\|\mathcal{A}\mathbf{e}\|}{\|\mathbf{e}\|} = \frac{\|\mathbf{r}\|}{\|\mathbf{e}\|} \leq \lambda_{\max} \Rightarrow \lambda_{\max}^{-1} \|\mathbf{r}\| \leq \|\mathbf{e}\| \leq \lambda_{\min}^{-1} \|\mathbf{r}\|.$$

Since \mathbf{u} is a solution, we have

$$\lambda_{\min} \leq \frac{\|\mathcal{A}\mathbf{u}\|}{\|\mathbf{u}\|} = \frac{\|\mathbf{f}\|}{\|\mathbf{u}\|} \leq \lambda_{\max} \Rightarrow \lambda_{\max}^{-1} \|\mathbf{f}\| \leq \|\mathbf{u}\| \leq \lambda_{\min}^{-1} \|\mathbf{f}\|$$

Thus,

$$\frac{\lambda_{\min}}{\lambda_{\max}} \frac{\|\mathbf{r}\|}{\|\mathbf{f}\|} \leq \frac{\|\mathbf{e}\|}{\|\mathbf{u}\|} \leq \frac{\lambda_{\max}}{\lambda_{\min}} \frac{\|\mathbf{r}\|}{\|\mathbf{f}\|}.$$

Key "residual–error" relation

Since

$$\frac{\lambda_{\max}}{\lambda_{\min}} = \mathbf{Cond} \mathcal{A},$$

we arrive at the basic relation where the matrix condition number serves as an important factor

$$\boxed{(\mathbf{Cond} \mathcal{A})^{-1} \frac{\|\mathbf{r}\|}{\|\mathbf{f}\|} \leq \frac{\|\mathbf{e}\|}{\|\mathbf{u}\|} \leq \mathbf{Cond} \mathcal{A} \frac{\|\mathbf{r}\|}{\|\mathbf{f}\|}.} \quad (43)$$

Thus, the relative error is controlled by the relative value of the residual. However, the bounds deteriorates when the conditional number is large.

In principle, the above consideration can be extended to a wider set of linear problems, where

$$\mathcal{A} \in \mathcal{L}(\mathbf{X}, \mathbf{Y})$$

is a coercive linear operator acting from a Banach space \mathbf{X} to another space \mathbf{Y} and \mathbf{f} is a given element of \mathbf{Y} .

However, if \mathcal{A} is related to a boundary-value problem, then one should properly define the spaces \mathbf{X} and \mathbf{Y} and find a practically meaningful analog of the estimate (43).

Elliptic equations

Let $\mathcal{A} : \mathbf{X} \rightarrow \mathbf{Y}$ be a linear elliptic operator. Consider the boundary-value problem

$$\mathcal{A}\mathbf{u} + \mathbf{f} = \mathbf{0} \quad \text{in } \Omega, \quad \mathbf{u} = \mathbf{u}_0 \quad \text{on } \partial\Omega.$$

Assume that $\mathbf{v} \in \mathbf{X}$ is an approximation of \mathbf{u} . Then, we should measure the error in \mathbf{X} and the residual in \mathbf{Y} , so that the principal form of the estimate is

$$\|\mathbf{v} - \mathbf{u}\|_{\mathbf{X}} \leq \mathbf{C} \|\mathcal{A}\mathbf{v} + \mathbf{f}\|_{\mathbf{Y}}, \quad (44)$$

where the constant \mathbf{C} is independent of \mathbf{v} . The key question is as follows:

Which spaces \mathbf{X} and \mathbf{Y} should we choose for a particular boundary-value problem?

Consider the problem

$$\Delta \mathbf{u} + \mathbf{f} = \mathbf{0} \quad \text{in } \Omega, \quad \mathbf{u} = \mathbf{0} \quad \text{on } \partial\Omega,$$

with $\mathbf{f} \in \mathbf{L}^2(\Omega)$. The generalized solution satisfies the relation

$$\int_{\Omega} \nabla \mathbf{u} \cdot \nabla \mathbf{w} \, d\mathbf{x} = \int_{\Omega} \mathbf{f} \mathbf{w} \, d\mathbf{x} \quad \forall \mathbf{w} \in \mathbf{V}_0 := \mathring{\mathbf{H}}^1(\Omega),$$

which implies the **energy estimate**

$$\|\nabla \mathbf{u}\|_{2,\Omega} \leq \mathbf{C}_{\Omega} \|\mathbf{f}\|_{2,\Omega}.$$

Here \mathbf{C}_{Ω} is a constant in the Friedrichs-Steklov inequality. Assume that an approximation $\mathbf{v} \in \mathbf{V}_0$ and $\Delta \mathbf{v} \in \mathbf{L}^2(\Omega)$. Then,

$$\int_{\Omega} \nabla(\mathbf{u} - \mathbf{v}) \cdot \nabla \mathbf{w} \, d\mathbf{x} = \int_{\Omega} (\mathbf{f} + \Delta \mathbf{v}) \mathbf{w} \, d\mathbf{x}, \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

Setting $\mathbf{w} = \mathbf{u} - \mathbf{v}$, we obtain the estimate

$$\|\nabla(\mathbf{u} - \mathbf{v})\|_{2,\Omega} \leq \mathbf{C}_\Omega \|\mathbf{f} + \Delta \mathbf{v}\|_{2,\Omega}, \quad (45)$$

whose right-hand side of (45) is formed by the \mathbf{L}^2 -norm of the residual. However, usually a sequence of approximations $\{\mathbf{v}_k\}$ converges to \mathbf{u} only in the energy space, i.e.,

$$\{\mathbf{v}_k\} \rightarrow \mathbf{u} \quad \text{in } \mathbf{H}^1(\Omega),$$

so that $\|\Delta \mathbf{v}_k + \mathbf{f}\|$ may not converge to zero !

This means that the **consistency** (the key property of any practically meaningful estimate) is lost.

Which norm of the residual leads to a consistent estimate of the error in the energy norm?

To find it, we should consider Δ not as $\mathbf{H}^2 \rightarrow \mathbf{L}^2$ mapping, but as $\mathbf{H}^1 \rightarrow \mathbf{H}^{-1}$ mapping. For this purpose we use the integral identity

$$\int_{\Omega} \nabla \mathbf{u} \cdot \nabla \mathbf{w} \, d\mathbf{x} = \langle \mathbf{f}, \mathbf{w} \rangle, \quad \forall \mathbf{w} \in \mathbf{V}_0 := \mathring{\mathbf{H}}^1(\Omega).$$

Here, $\nabla \mathbf{u} \in \mathbf{L}^2$, so that it has derivatives in \mathbf{H}^{-1} and we consider the above as equivalence of two distributions on all trial functions $\mathbf{w} \in \mathbf{V}_0$.

By $\langle \mathbf{f}, \mathbf{w} \rangle \leq \mathbf{I} \mathbf{f} \mathbf{I} \|\nabla \mathbf{w}\|_{2,\Omega}$, we obtain another "energy estimate"

$$\|\nabla \mathbf{u}\|_{2,\Omega} \leq \mathbf{I} \mathbf{f} \mathbf{I}.$$

Consistent residual estimate

Let $\mathbf{v} \in \mathbf{V}_0$ be an approximation of \mathbf{u} . We have

$$\begin{aligned} \int_{\Omega} \nabla(\mathbf{u} - \mathbf{v}) \cdot \nabla \mathbf{w} \, dx &= \int_{\Omega} (\mathbf{f}\mathbf{w} - \nabla \mathbf{v} \cdot \nabla \mathbf{w}) \, dx = \\ &= \langle \Delta \mathbf{v} + \mathbf{f}, \mathbf{w} \rangle, \quad \mathbf{f} + \Delta \mathbf{v} \in \mathbf{H}^{-1}(\Omega). \end{aligned}$$

By setting $\mathbf{w} = \mathbf{v} - \mathbf{u}$, we obtain

$$\|\nabla(\mathbf{u} - \mathbf{v})\|_{2,\Omega} \leq \mathbf{I} \mathbf{f} + \Delta \mathbf{v} \mathbf{I}. \quad (46)$$

where

$$\begin{aligned} \mathbf{I} \mathbf{f} + \Delta \mathbf{v} \mathbf{I} &= \sup_{\varphi \in \mathring{\mathbf{H}}^1(\Omega)} \frac{|\langle \mathbf{f} + \Delta \mathbf{v}, \varphi \rangle|}{\|\nabla \varphi\|} = \\ &= \sup_{\varphi \in \mathring{\mathbf{H}}^1(\Omega)} \frac{|\int_{\Omega} \nabla(\mathbf{u} - \mathbf{v}) \cdot \nabla \varphi|}{\|\nabla \varphi\|} \leq \sup_{\varphi \in \mathring{\mathbf{H}}^1(\Omega)} \frac{\|\nabla(\mathbf{u} - \mathbf{v})\| \|\nabla \varphi\|}{\|\nabla \varphi\|} \leq \|\nabla(\mathbf{u} - \mathbf{v})\| \end{aligned}$$

Thus, for the problem considered

$$\|\nabla(\mathbf{u} - \mathbf{v})\|_{2,\Omega} = \mathbf{I} \mathbf{f} + \Delta \mathbf{v} \mathbf{I} !!! \quad (47)$$

From (47), it readily follows that

$$\mathbf{I} \mathbf{f} + \Delta \mathbf{v}_k \mathbf{I} \rightarrow \mathbf{0} \quad \text{as} \quad \{\mathbf{v}_k\} \rightarrow \mathbf{u} \text{ in } \mathbf{H}^1.$$

We observe that the estimate (47) is **consistent**.

Diffusion equation

Similar estimates can be derived for

$$\mathcal{A}\mathbf{u} + \mathbf{f} = \mathbf{0}, \quad \text{in } \Omega, \quad \mathbf{u} = \mathbf{0} \text{ on } \partial\Omega,$$

where

$$\mathcal{A}\mathbf{u} = \mathbf{div} \mathbf{A} \nabla \mathbf{u} := \sum_{i,j=1}^d \frac{\partial}{\partial \mathbf{x}_i} \left(\mathbf{a}_{ij}(\mathbf{x}) \frac{\partial \mathbf{u}}{\partial \mathbf{x}_j} \right),$$

$$\mathbf{a}_{ij}(\mathbf{x}) = \mathbf{a}_{ji}(\mathbf{x}) \in \mathbf{L}^\infty(\Omega),$$

$$\lambda_{\min} |\boldsymbol{\eta}|^2 \leq \mathbf{a}_{ij}(\mathbf{x}) \eta_i \eta_j \leq \lambda_{\max} |\boldsymbol{\eta}|^2, \quad \forall \boldsymbol{\eta} \in \mathbb{R}^n, \mathbf{x} \in \Omega,$$

$$\lambda_{\max} \geq \lambda_{\min} \geq 0.$$

Let $\mathbf{v} \in \mathbf{V}_0$ be an approximation of \mathbf{u} . Then,

$$\int_{\Omega} \mathbf{A} \nabla(\mathbf{u} - \mathbf{v}) \cdot \nabla \mathbf{w} \, dx = \int_{\Omega} (\mathbf{f} \mathbf{w} - \mathbf{A} \nabla \mathbf{v} \cdot \nabla \mathbf{w}) \, dx, \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

Again, the right-hand side of this relation is a bounded linear functional on \mathbf{V}_0 , i.e.,

$$\mathbf{f} + \mathbf{div}(\mathbf{A} \nabla \mathbf{v}) \in \mathbf{H}^{-1}.$$

Hence, we have the relation

$$\int_{\Omega} \mathbf{A} \nabla(\mathbf{u} - \mathbf{v}) \cdot \nabla \mathbf{w} \, dx = \langle \mathbf{f} + \mathbf{div}(\mathbf{A} \nabla \mathbf{v}), \mathbf{w} \rangle, \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

Setting $\mathbf{w} = \mathbf{u} - \mathbf{v}$, we derive the estimate

$$\|\nabla(\mathbf{u} - \mathbf{v})\|_{2,\Omega} \leq \lambda_{\min}^{-1} \|\mathbf{f} + \mathbf{div}(\mathbf{A} \nabla \mathbf{v})\|. \quad (48)$$

Next,

$$\begin{aligned} \mathbf{I} \mathbf{f} + \mathbf{div}(\mathbf{A}\nabla\mathbf{v}) \mathbf{I} &= \sup_{\varphi \in \mathring{\mathbf{H}}^1(\Omega)} \frac{|\langle \mathbf{f} + \mathbf{div}(\mathbf{A}\nabla\mathbf{v}), \varphi \rangle|}{\|\nabla\varphi\|_{2,\Omega}} = \\ &= \sup_{\varphi \in \mathring{\mathbf{H}}^1(\Omega)} \frac{|\int_{\Omega} \mathbf{A}\nabla(\mathbf{u} - \mathbf{v}) \cdot \nabla\varphi \, \mathbf{d}\mathbf{x}|}{\|\nabla\varphi\|_{2,\Omega}} \leq \lambda_{\max} \|\nabla(\mathbf{u} - \mathbf{v})\|_{2,\Omega}. \end{aligned} \quad (49)$$

Combining (48) and (49) we obtain

$$\boxed{\lambda_{\max}^{-1} \mathbf{I} \mathbf{R}(\mathbf{v}) \mathbf{I} \leq \|\nabla(\mathbf{u} - \mathbf{v})\|_{2,\Omega} \leq \lambda_{\min}^{-1} \mathbf{I} \mathbf{R}(\mathbf{v}) \mathbf{I}}, \quad (50)$$

where $\mathbf{R}(\mathbf{v}) = \mathbf{f} + \mathbf{div}(\mathbf{A}\nabla\mathbf{v}) \in \mathbf{H}^{-1}(\Omega)$. We see that upper and lower bounds of the error can be evaluated in terms of the negative norm of $\mathbf{R}(\mathbf{v})$.

Main goal

We observe that to find guaranteed bounds of the error reliable estimates of $\mathbf{R}(\mathbf{v})$ are required.

In essence, a posteriori error estimates derived in 70-90' for Finite Element Methods (FEM) offer several approaches to the evaluation of $\mathbf{R}(\mathbf{v})$. We consider them starting with the so-called **explicit residual method** where such estimates are obtained with help of two key points:

- Galerkin orthogonality property;
- $\mathbf{H}^1 \rightarrow \mathbf{V}_h$ interpolation estimates by Clément.

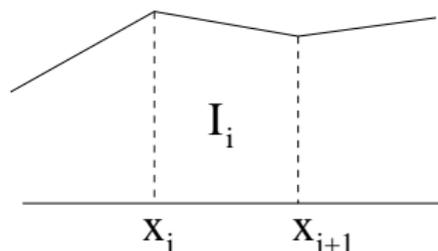
Explicit residual method in 1D case

Take the simplest model

$$(\alpha \mathbf{u}')' + \mathbf{f} = \mathbf{0}, \quad \mathbf{u}(\mathbf{0}) = \mathbf{u}(\mathbf{1}).$$

Let $\mathbf{I} := (\mathbf{0}, \mathbf{1})$, $\mathbf{f} \in \mathbf{L}^2(\mathbf{I})$, $\alpha(\mathbf{x}) \in \mathbf{C}(\bar{\mathbf{I}}) \geq \alpha_0 > \mathbf{0}$. Divide \mathbf{I} into a number of subintervals $\mathbf{I}_i = (\mathbf{x}_i, \mathbf{x}_{i+1})$, where $\mathbf{x}_0 = \mathbf{0}$, $\mathbf{x}_{N+1} = \mathbf{1}$, and $|\mathbf{x}_{i+1} - \mathbf{x}_i| = \mathbf{h}_i$.

Assume that $\mathbf{v} \in \mathring{\mathbf{H}}^1(\mathbf{I})$ and it is smooth on any interval \mathbf{I}_j .



In this case,

$$\begin{aligned}
 \mathbf{I} \mathbf{R}(\mathbf{v}) \mathbf{I} &= \sup_{\mathbf{w} \in \mathbf{V}_0(I), \mathbf{w} \neq \mathbf{0}} \frac{\int_0^1 (-\alpha \mathbf{v}' \mathbf{w}' + \mathbf{f} \mathbf{w}) \mathbf{d}x}{\|\mathbf{w}'\|_{2,I}} = \\
 &= \sup_{\mathbf{w} \in \overset{\circ}{\mathbf{H}}^1(I); \mathbf{w} \neq \mathbf{0}} \frac{\sum_{i=0}^N \int_{I_i} (-\alpha \mathbf{v}' \mathbf{w}' + \mathbf{f} \mathbf{w}) \mathbf{d}x}{\|\mathbf{w}'\|_{2,I}} = \\
 &= \sup_{\mathbf{w} \in \mathbf{V}_0(I), \mathbf{w} \neq \mathbf{0}} \frac{\sum_{i=0}^N \int_{I_i} \mathbf{r}_i(\mathbf{v}) \mathbf{w} \mathbf{d}x + \sum_{i=1}^N \alpha(\mathbf{x}_i) \mathbf{w}(\mathbf{x}_i) \mathbf{j}(\mathbf{v}'(\mathbf{x}_i))}{\|\mathbf{w}'\|_{2,I}},
 \end{aligned}$$

where $\mathbf{j}(\phi(\mathbf{x})) := \phi(\mathbf{x} + \mathbf{0}) - \phi(\mathbf{x} - \mathbf{0})$ is the "jump-function" and $\mathbf{r}_i(\mathbf{v}) = (\alpha \mathbf{v}')' + \mathbf{f}$ is the residual on I_i .

For arbitrary \mathbf{v} we can hardly get an upper bound for this supremum.

Use Galerkin orthogonality

Assume that $\mathbf{v} = \mathbf{u}_h$, i.e., it is the *Galerkin approximation* obtained on a finite-dimensional subspace \mathbf{V}_{0h} formed by piecewise polynomial continuous functions. Since

$$\int_I \alpha \mathbf{u}'_h \mathbf{w}'_h \, dx - \int_I \mathbf{f} \mathbf{w}_h \, dx = \mathbf{0} \quad \forall \mathbf{w}_h \in \mathbf{V}_{0h}.$$

we may add the left-hand side with any w_h to the numerator what gives

$$\mathbf{I} \mathbf{R}(\mathbf{u}_h) \mathbf{I} = \sup_{\mathbf{w} \in \mathbf{V}_0(I)} \frac{\int_0^1 (-\alpha \mathbf{u}'_h(\mathbf{w} - \boldsymbol{\pi}_h \mathbf{w})' + \mathbf{f}(\mathbf{w} - \boldsymbol{\pi}_h \mathbf{w})) \, dx}{\|\mathbf{w}'\|_{2,I}},$$

where $\boldsymbol{\pi}_h : \mathbf{V}_0 \rightarrow \mathbf{V}_{0h}$ is the interpolation operator defined by the conditions $\boldsymbol{\pi}_h \mathbf{v} \in \mathbf{V}_{0h}$, $\boldsymbol{\pi}_h \mathbf{v}(0) = \boldsymbol{\pi}_h \mathbf{v}(1) = \mathbf{0}$ and

$$\boldsymbol{\pi}_h \mathbf{v}(\mathbf{x}_i) = \mathbf{v}(\mathbf{x}_i), \quad \forall \mathbf{x}_i, \quad \mathbf{i} = 1, 2, \dots, \mathbf{N}.$$

Integrating by parts

Now, we have

$$\mathbf{I} \mathbf{R}(\mathbf{u}_h) \mathbf{I} = \sup_{\mathbf{w} \in \mathbf{V}_0(\Omega)} \left\{ \frac{\sum_{i=0}^N \int_{I_i} \mathbf{r}_i(\mathbf{u}_h)(\mathbf{w} - \pi_h \mathbf{w}) \, d\mathbf{x}}{\|\mathbf{w}'\|_{2,1}} + \frac{\sum_{i=1}^N \alpha(\mathbf{x}_i)(\mathbf{w}(\mathbf{x}_i) - \pi_h \mathbf{w}(\mathbf{x}_i)) \mathbf{j}(\mathbf{u}'_h(\mathbf{x}_i))}{\|\mathbf{w}'\|_{2,1}} \right\}.$$

Since $\mathbf{w}(\mathbf{x}_i) - \pi_h \mathbf{w}(\mathbf{x}_i) = \mathbf{0}$, the second sum vanishes. For first one we have

$$\sum_{i=0}^N \int_{I_i} \mathbf{r}_i(\mathbf{u}_h)(\mathbf{w} - \pi_h \mathbf{w}) \, d\mathbf{x} \leq \sum_{i=0}^N \|\mathbf{r}_i(\mathbf{u}_h)\|_{2,I_i} \|\mathbf{w} - \pi_h \mathbf{w}\|_{2,I_i}.$$

Since for $\mathbf{w} \in \mathring{\mathbf{H}}^1(\mathbf{I}_i)$

$$\|\mathbf{w} - \pi_h \mathbf{w}\|_{2,\mathbf{I}_i} \leq \mathbf{c}_i \|\mathbf{w}'\|_{2,\mathbf{I}_i},$$

we obtain for the numerator of the above quotient

$$\begin{aligned} \sum_{i=0}^N \int_{\mathbf{I}_i} \mathbf{r}_i(\mathbf{u}_h)(\mathbf{w} - \pi_h \mathbf{w}) \, d\mathbf{x} &\leq \sum_{i=0}^N \mathbf{c}_i \|\mathbf{r}_i(\mathbf{u}_h)\|_{2,\mathbf{I}_i} \|\mathbf{w}'\|_{2,\mathbf{I}_i} \leq \\ &\leq \left(\sum_{i=0}^N \mathbf{c}_i^2 \|\mathbf{r}_i(\mathbf{u}_h)\|_{2,\mathbf{I}_i}^2 \right)^{1/2} \|\mathbf{w}'\|_{2,\mathbf{I}}, \end{aligned}$$

which implies the desired upper bound

$$\boxed{\mathbf{I} \mathbf{R}(\mathbf{u}_h) \mathbf{I} \leq \left(\sum_{i=0}^N \mathbf{c}_i^2 \|\mathbf{r}_i(\mathbf{u}_h)\|_{2,\mathbf{I}_i}^2 \right)^{1/2}}. \quad (51)$$

This bound is the sum of local residuals $\mathbf{r}_i(\mathbf{u}_h)$ with weights given by the **interpolation constants** \mathbf{c}_i .

Interpolation constants

For piecewise affine approximations, the interpolation constants \mathbf{c}_i are easy to find. Indeed, let γ_i be a constant that satisfies the condition

$$\inf_{\mathbf{w} \in \mathring{\mathbf{H}}^1(I_i)} \frac{\|\mathbf{w}'\|_{2,I_i}^2}{\|\mathbf{w} - \pi_h \mathbf{w}\|_{2,I_i}^2} \geq \gamma_i.$$

Then, for all $\mathbf{w} \in \mathring{\mathbf{H}}^1(I_i)$, we have

$$\|\mathbf{w} - \pi_h \mathbf{w}\|_{2,I_i} \leq \gamma_i^{-1/2} \|\mathbf{w}'\|_{2,I_i}$$

and one can set $\mathbf{c}_i = \gamma_i^{-1/2}$.

Let us estimate γ_{I_i} .

Note that

$$\int_{x_i}^{x_{i+1}} |\mathbf{w}'|^2 \, d\mathbf{x} = \int_{x_i}^{x_{i+1}} |(\mathbf{w} - \pi_h \mathbf{w})' + (\pi_h \mathbf{w})'|^2 \, d\mathbf{x},$$

where $(\pi_h \mathbf{w})'$ is constant on (x_i, x_{i+1}) . Therefore,

$$\int_{x_i}^{x_{i+1}} (\mathbf{w} - \pi_h \mathbf{w})' (\pi_h \mathbf{w})' \, d\mathbf{x} = 0$$

and

$$\begin{aligned} \int_{x_i}^{x_{i+1}} |\mathbf{w}'|^2 \, d\mathbf{x} &= \int_{x_i}^{x_{i+1}} |(\mathbf{w} - \pi_h \mathbf{w})'|^2 \, d\mathbf{x} + \int_{x_i}^{x_{i+1}} |(\pi_h \mathbf{w})'|^2 \, d\mathbf{x} \geq \\ &\geq \int_{x_i}^{x_{i+1}} |(\mathbf{w} - \pi_h \mathbf{w})'|^2 \, d\mathbf{x}. \end{aligned}$$

Interpolation constants in 1D problem

Thus, we have

$$\begin{aligned} \inf_{\mathbf{w} \in \mathring{H}^1(I_i)} \frac{\int_{x_i}^{x_{i+1}} |\mathbf{w}'|^2 \, dx}{\int_{x_i}^{x_{i+1}} |\mathbf{w} - \pi_h \mathbf{w}|^2 \, dx} &\geq \inf_{\mathbf{w} \in \mathring{H}^1(I_i)} \frac{\int_{x_i}^{x_{i+1}} |(\mathbf{w} - \pi_h \mathbf{w})'|^2 \, dx}{\int_{x_i}^{x_{i+1}} |\mathbf{w} - \pi_h \mathbf{w}|^2 \, dx} \geq \\ &\geq \inf_{\eta \in \mathring{H}^1(I_i)} \frac{\int_{x_i}^{x_{i+1}} |\eta'|^2 \, dx}{\int_{x_i}^{x_{i+1}} |\eta|^2 \, dx} = \frac{\pi^2}{h_i^2}, \end{aligned}$$

so that $\gamma_i = \pi^2/h_i^2$ and $\mathbf{c}_i = h_i/\pi$.

Remark. To prove the very last relation we note that

$$\inf_{\eta \in \mathring{H}^1((0,h))} \frac{\int_0^h |\eta'|^2 \, dx}{\int_0^h |\eta|^2 \, dx} = \frac{\pi^2}{h^2}$$

is attained on the eigenfunction $\sin \frac{\pi}{h} x$, of the problem $\phi'' + \lambda \phi = 0$ on $(0, h)$.

Task 4

Solve a boundary-value problem

$$\begin{aligned}(\alpha \mathbf{v}')' &= \mathbf{f}, \\ \mathbf{v}(\mathbf{0}) &= \mathbf{a}, \quad \mathbf{v}(\mathbf{1}) = \mathbf{b}\end{aligned}$$

with certain $\alpha(x) > 0$, \mathbf{f} , \mathbf{a} , and \mathbf{b} by the finite element method with uniform elements (i.e., $\mathbf{h} = \mathbf{1}/\mathbf{N}$). Apply the residual method and compare the errors computed with the true error distribution.

Residual method in 2D case

Let Ω be represented as a union \mathcal{T}_h of simplexes \mathbf{T}_i . For the sake of simplicity, assume that $\bar{\Omega} = \cup_{i=1}^N \bar{\mathbf{T}}_i$ and \mathbf{V}_{0h} consists of piecewise affine continuous functions. Then the Galerkin approximation \mathbf{u}_h satisfies the relation

$$\int_{\Omega} \mathbf{A} \nabla \mathbf{u}_h \cdot \nabla \mathbf{w}_h \, d\mathbf{x} = \int_{\Omega} \mathbf{f} \mathbf{w}_h \, d\mathbf{x}, \quad \forall \mathbf{w}_h \in \mathbf{V}_{0h},$$

where

$$\mathbf{V}_{0h} = \{ \mathbf{w}_h \in \mathbf{V}_0 \mid \mathbf{w}_h \in \mathbf{P}^1(\mathbf{T}_i), \mathbf{T}_i \in \mathcal{F}_h \}.$$

In this case, negative norm of the residual is

$$\| \mathbf{R}(\mathbf{u}_h) \| = \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\int_{\Omega} (\mathbf{f}\mathbf{w} - \mathbf{A}\nabla\mathbf{u}_h \cdot \nabla\mathbf{w}) \, d\mathbf{x}}{\|\nabla\mathbf{w}\|_{2,\Omega}}.$$

Let $\pi : \overset{\circ}{H}^1 \rightarrow V_{0h}$ be a continuous interpolation operator. Then, for the **Galerkin approximation**

$$\| \mathbf{R}(\mathbf{u}_h) \| = \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\int_{\Omega} (\mathbf{f}(\mathbf{w} - \pi_h\mathbf{w}) - \mathbf{A}\nabla\mathbf{u}_h \cdot \nabla(\mathbf{w} - \pi_h\mathbf{w})) \, d\mathbf{x}}{\|\nabla\mathbf{w}\|_{2,\Omega}}.$$

For finite element approximations such a type projection operators has been constructed. One of the most known was suggested in

[Ph. Clément. Approximations by finite element functions using local regularization, *RAIRO Anal. Numér.*, 9\(1975\).](#)

and is often called the **Clement's interpolation operator**. Its properties play an important role in the a posteriori error estimation method considered.

Clement's Interpolation operator

Let \mathbf{E}_{ij} denote the common edge of the simplexes \mathbf{T}_i and \mathbf{T}_j . If \mathbf{s} is an inner node of the triangulation \mathcal{F}_h , then ω_s denotes the set of all simplexes having this node.

For any \mathbf{s} , we find a polynomial $\mathbf{p}_s(\mathbf{x}) \in \mathbf{P}^1(\omega_s)$ such that

$$\int_{\omega_s} (\mathbf{v} - \mathbf{p}_s) \mathbf{q} \, d\mathbf{x} = \mathbf{0} \quad \forall \mathbf{q} \in \mathbf{P}^1(\omega_s).$$

Now, the interpolation operator π_h is defined by setting

$$\pi_h \mathbf{v}(\mathbf{x}_s) = \mathbf{p}(\mathbf{x}_s), \quad \forall \mathbf{x}_s \in \Omega,$$

$$\pi_h \mathbf{v}(\mathbf{x}_s) = \mathbf{0}, \quad \forall \mathbf{x}_s \in \partial\Omega.$$

It is a linear and continuous mapping of $\mathring{\mathbf{H}}^1(\Omega)$ to the space of piecewise affine continuous functions.

Interpolation estimates in 2D

Moreover, it is subject to the relations

$$\|\mathbf{v} - \pi_{\mathbf{h}}\mathbf{v}\|_{2, \mathbf{T}_i} \leq \mathbf{c}_i^T \text{diam}(\mathbf{T}_i) \|\mathbf{v}\|_{1,2, \omega_{\mathbf{N}}(\mathbf{T}_i)}, \quad (52)$$

$$\|\mathbf{v} - \pi_{\mathbf{h}}\mathbf{v}\|_{2, \mathbf{E}_{ij}} \leq \mathbf{c}_{ij}^E |\mathbf{E}_{ij}|^{1/2} \|\mathbf{v}\|_{1,2, \omega_{\mathbf{E}}(\mathbf{T}_i)}, \quad (53)$$

where $\omega_{\mathbf{N}}(\mathbf{T}_i)$ is the union of all simplexes having at least *one common node* with \mathbf{T}_i and $\omega_{\mathbf{E}}(\mathbf{T}_i)$ is the union of all simplexes having *a common edge* with \mathbf{T}_i .

Interpolation constants \mathbf{c}_i^T and \mathbf{c}_{ij}^E are LOCAL and depend on the shape of patches $\omega_{\mathbf{N}}(\mathbf{T}_i)$ and $\omega_{\mathbf{E}}(\mathbf{T}_i)$.

Quotient relations for the constants

Evaluation of \mathbf{c}_i^T and \mathbf{c}_{ij}^E requires finding *exact lower bounds* of the following variational problems:

$$\gamma_i^T := \inf_{\mathbf{w} \in \mathbf{V}_0} \frac{\|\mathbf{w}\|_{1,2,\omega_N(\mathbf{T}_i)}}{\|\mathbf{w} - \pi_h \mathbf{w}\|_{2,\mathbf{T}_i}} \text{diam}(\mathbf{T}_i)$$

and

$$\gamma_{ij}^E := \inf_{\mathbf{w} \in \mathbf{V}_0} \frac{\|\mathbf{w}\|_{1,2,\omega_E(\mathbf{T}_i)}}{\|\mathbf{w} - \pi_h \mathbf{w}\|_{2,\mathbf{E}_{ij}}} |\mathbf{E}_{ij}|^{1/2}.$$

Certainly, we can replace \mathbf{V}_0 by $\mathbf{H}^1(\omega_N(\mathbf{T}_i))$ and $\mathbf{H}^1(\omega_E(\mathbf{T}_i))$, respectively, but, anyway finding the constants amounts solving functional eigenvalue type problems !

Let $\boldsymbol{\sigma}_h = \mathbf{A}\nabla\mathbf{u}_h$. Then,

$$\mathbf{I} \mathbf{R}(\mathbf{u}_h) \mathbf{I} = \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\int_{\Omega} (\mathbf{f}(\mathbf{w} - \pi_h \mathbf{w}) - \boldsymbol{\sigma}_h \cdot \nabla(\mathbf{w} - \pi_h \mathbf{w})) \, \mathbf{d}\mathbf{x}}{\|\nabla \mathbf{w}\|_{2,\Omega}}.$$

If $\boldsymbol{\nu}_{ij}$ is the unit outward normal to \mathbf{E}_{ij} , then

$$\begin{aligned} \int_{\mathbf{T}_i} \boldsymbol{\sigma}_h \cdot \nabla(\mathbf{w} - \pi_h \mathbf{w}) \, \mathbf{d}\mathbf{x} &= \\ &= \sum_{\mathbf{E}_{ij} \subset \partial \mathbf{T}_i} \int_{\mathbf{E}_{ij}} (\boldsymbol{\sigma}_h \cdot \boldsymbol{\nu})(\mathbf{w} - \pi_h \mathbf{w}) \, \mathbf{d}\mathbf{s} - \int_{\mathbf{T}_i} \mathbf{div} \boldsymbol{\sigma}_h (\mathbf{w} - \pi_h \mathbf{w}) \, \mathbf{d}\mathbf{x}, \end{aligned}$$

Since on the boundary $\mathbf{w} - \pi_h \mathbf{w} = \mathbf{0}$, we obtain

$$\mathbf{I} \mathbf{R}(\mathbf{u}_h) \mathbf{I} = \sup_{\mathbf{w} \in \mathbf{V}_0} \left\{ \frac{\sum_{i=1}^N \int_{\mathbf{T}_i} (\mathbf{div} \boldsymbol{\sigma}_h + \mathbf{f})(\mathbf{w} - \pi_h \mathbf{w}) \, \mathbf{d}\mathbf{x}}{\|\nabla \mathbf{w}\|_{2,\Omega}} + \frac{\sum_{i=1}^N \sum_{j>i}^N \int_{\mathbf{E}_{ij}} \mathbf{j}(\boldsymbol{\sigma}_h \cdot \boldsymbol{\nu}_{ij})(\mathbf{w} - \pi_h \mathbf{w}) \, \mathbf{d}\mathbf{s}}{\|\nabla \mathbf{w}\|_{2,\Omega}} \right\}.$$

First term in sup

$$\begin{aligned} \int_{\mathbf{T}_i} (\mathbf{div} \boldsymbol{\sigma}_h + \mathbf{f})(\mathbf{w} - \pi_h \mathbf{w}) \, d\mathbf{x} &\leq \|\mathbf{div} \boldsymbol{\sigma}_h + \mathbf{f}\|_{2, \mathbf{T}_i} \|\mathbf{w} - \pi_h \mathbf{w}\|_{2, \mathbf{T}_i} \\ &\leq \mathbf{c}_i^T \|\mathbf{div} \boldsymbol{\sigma}_h + \mathbf{f}\|_{2, \mathbf{T}_i} \mathbf{diam}(\mathbf{T}_i) \|\mathbf{w}\|_{1,2, \omega_{\mathbf{N}}(\mathbf{T}_i)}, \end{aligned}$$

Then, the first sum is estimated as follows:

$$\begin{aligned} \sum_{i=1}^N \int_{\mathbf{T}_i} (\mathbf{div} \boldsymbol{\sigma}_h + \mathbf{f})(\mathbf{w} - \pi_h \mathbf{w}) \, d\mathbf{x} &\leq \\ &\leq \mathbf{d}_1 \left(\sum_{i=1}^N (\mathbf{c}_i^T)^2 \mathbf{diam}(\mathbf{T}_i)^2 \|\mathbf{div} \boldsymbol{\sigma}_h + \mathbf{f}\|_{2, \mathbf{T}_i}^2 \right)^{1/2} \|\mathbf{w}\|_{1,2, \Omega}, \end{aligned}$$

where the constant \mathbf{d}_1 depends on the maximal number of elements in the set $\omega_{\mathbf{N}}(\mathbf{T}_i)$.

Second term in sup

For the second one, we have

$$\begin{aligned}
 & \sum_{i=1}^N \sum_{j>i}^N \int_{\mathbf{E}_{ij}} \mathbf{j}(\boldsymbol{\sigma}_h \cdot \boldsymbol{\nu}_{ij})(\mathbf{w} - \pi_h \mathbf{w}) \, d\mathbf{x} \leq \\
 & \leq \sum_{i=1}^N \sum_{j>i}^N \|\mathbf{j}(\boldsymbol{\sigma}_h \cdot \boldsymbol{\nu}_{ij})\|_{2, \mathbf{E}_{ij}} \mathbf{c}_{ij}^E |\mathbf{E}_{ij}|^{1/2} \|\mathbf{w}\|_{1,2, \omega_{\mathbf{E}}(\mathbf{T}_i)} \leq \\
 & \leq \mathbf{d}_2 \left(\sum_{i=1}^N \sum_{j>i}^N \left(\mathbf{c}_{ij}^E \right)^2 |\mathbf{E}_{ij}| \|\mathbf{j}(\boldsymbol{\sigma}_h \cdot \boldsymbol{\nu}_{ij})\|_{2, \mathbf{E}_{ij}}^2 \right)^{1/2} \|\mathbf{w}\|_{1,2, \Omega},
 \end{aligned}$$

where \mathbf{d}_2 depends on the maximal number of elements in the set $\omega_{\mathbf{E}}(\mathbf{T}_i)$.

Residual type error estimate

By the above estimates we obtain

$$\begin{aligned} \mathbf{I} \mathbf{R}(\mathbf{u}_h) \mathbf{I} \leq \mathbf{C}_0 & \left(\left(\sum_{i=1}^N (\mathbf{c}_i^T)^2 \text{diam}(\mathbf{T}_i)^2 \|\mathbf{div} \boldsymbol{\sigma}_h + \mathbf{f}\|_{2, \mathbf{T}_i}^2 \right)^{1/2} + \right. \\ & \left. + \left(\sum_{i=1}^N \sum_{j>i}^N (\mathbf{c}_{ij}^E)^2 |\mathbf{E}_{ij}| \|\mathbf{j}(\boldsymbol{\sigma}_h \cdot \boldsymbol{\nu}_{ij})\|_{2, \mathbf{E}_{ij}}^2 \right)^{1/2} \right). \quad (54) \end{aligned}$$

Here $\mathbf{C}_0 = \mathbf{C}_0(\mathbf{d}_1, \mathbf{d}_2)$. We observe that the right-hand side is the sum of local quantities (usually denoted by $\boldsymbol{\eta}(\mathbf{T}_i)$) multiplied by constants depending on properties of the chosen splitting \mathcal{F}_h .

Error indicator for quasi-uniform meshes

For quasi-uniform meshes all generic constants \mathbf{c}_i^T have approximately the same value and can be replaced by a single constant \mathbf{c}_1 . If the constants \mathbf{c}_{ij}^E are also estimated by a single constant \mathbf{c}_2 , then we have

$$\mathbf{I} \mathbf{R}(\mathbf{u}_h) \mathbf{I} \leq \mathbf{C} \left(\sum_{i=1}^N \eta^2(\mathbf{T}_i) \right)^{1/2}, \quad (55)$$

where $\mathbf{C} = \mathbf{C}(\mathbf{c}_1, \mathbf{c}_2, \mathbf{C}_0)$ and

$$\eta^2(\mathbf{T}_i) = \mathbf{c}_1^2 \text{diam}(\mathbf{T}_i)^2 \|\text{div} \sigma_h + \mathbf{f}\|_{2, \mathbf{T}_i}^2 + \frac{\mathbf{c}_2^2}{2} \sum_{\mathbf{E}_{ij} \subset \partial \mathbf{T}_i} |\mathbf{E}_{ij}| \|\mathbf{j}(\sigma_h \cdot \nu_{ij})\|_{2, \mathbf{E}_{ij}}^2.$$

The multiplier 1/2 arises, because any interior edge is common for two elements.

Comment 1

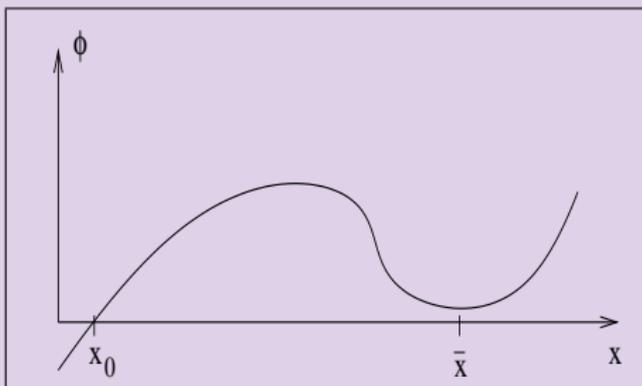
General form of the residual type a posteriori error estimates is as follows:

$$\|\mathbf{u} - \mathbf{u}_h\| \leq \mathbf{M}(\mathbf{u}_k, \mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_N, \mathcal{D}),$$

where \mathcal{D} is the data set, \mathbf{u}_h is the **Galerkin approximation**, and $\mathbf{c}_i, i = 1, 2, \dots, N$ are the **interpolation constants**. The constants depend on the **mesh** and properties of the special type interpolation operator. The number N depends on the dimension of \mathbf{V}_h and may be rather large. If the constants are not sharply defined, then this functional is not more than a certain error indicator. However, in many cases it successfully works and was used in numerous researches.

Comment 2

It is worth noting that for nonlinear problems the dependence between the error and the respective residual is much more complicated. A simple example below shows that the value of the residual may fail to control the distance to the exact solution.



References

It is commonly accepted that this approach brings its origin from the papers

I. Babuska and W. C. Rheinboldt. A-posteriori error estimates for the finite element method. *Internat. J. Numer. Meth. Engrg.*, 12(1978).

I. Babuska and W. C. Rheinboldt. Error estimates for adaptive finite element computations. *SIAM J.Numer. Anal.*, 15(1978). Detailed mathematical analysis of this error estimation method can be found in R. Verfürth. *A review of a posteriori error estimation and adaptive mesh-refinement techniques* Wiley and Sons, Teubner, New-York, 1996.

Also, we recommend the books

M. Ainsworth and T. Oden. *A posteriori error estimation in finite element analysis*, Wiley and Sons, New York, 2000.

K. Eriksson, D. Estep, P. Hansbo, and C. Johnson, *Computational differential equations*, Cambridge University Press, Cambridge, 1996

I. Babuska and T. Strouboulis, *The finite element method and its reliability*, Oxford University Press, New York, 2001.

A posteriori methods based on post-processing

Post-processing of approximate solutions is a numerical procedure intended to modify already computed solution in such a way that the post-processed function would fit some **a priori known properties** much better than the original one.

Preliminaries

Let \mathbf{e} denotes the **error** of an approximate solution $\mathbf{v} \in \mathbf{V}$ and $\mathcal{E}(\mathbf{v}) : \mathbf{V} \rightarrow \mathbf{R}_+$ denotes the value of an **error estimator** computed on \mathbf{v} .

Definition

The estimator is said to be **equivalent to the error** for the approximations \mathbf{v} from a certain subset $\tilde{\mathbf{V}}$ if

$$c_1 \mathcal{E}(\mathbf{v}) \leq \|\mathbf{e}\| \leq c_2 \mathcal{E}(\mathbf{v}) \quad \forall \mathbf{v} \in \tilde{\mathbf{V}}$$

Definition

The ratio

$$\mathbf{i}_{\text{eff}} := \mathbf{1} + \frac{\mathcal{E}(\mathbf{v}) - \|\mathbf{e}\|}{\|\mathbf{e}\|}$$

is called the **effectivity index** of the estimator \mathcal{E} .

Ideal estimator has $\mathbf{i}_{\text{eff}} = \mathbf{1}$. However, in real life situations it is hardly possible, so that values \mathbf{i}_{eff} in the diapason from 1 to 2-3 are considered as quite good.

In FEM methods with mesh size \mathbf{h} one other term is often used:

Definition

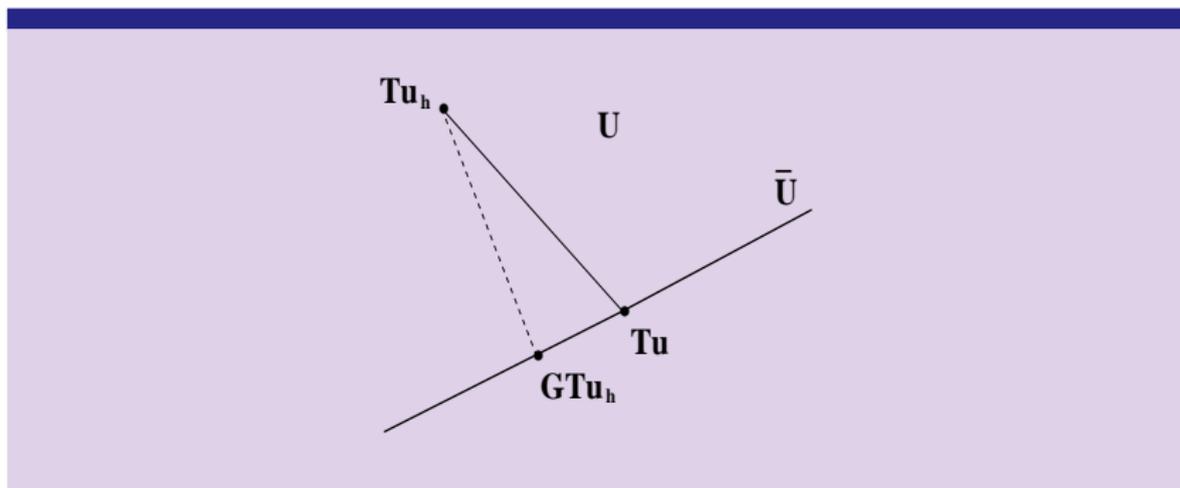
The estimator \mathcal{E} is called **asymptotically equivalent to the error** if for a sequence of approximate solutions $\{u_h\}$ obtained on consequently refined meshes there holds the relation

$$\inf_{\mathbf{h} \rightarrow 0} \frac{\mathcal{E}(\mathbf{u}_h)}{\|\mathbf{u} - \mathbf{u}_h\|} = \mathbf{1}$$

It is clear that an estimator may be asymptotically exact for one sequence of approximate solutions (e.g. computed on regular meshes) and not exact for another one.

General outlook

Typically, the function $T\mathbf{u}_h$ (where T is a certain linear operator, e.g., ∇) lies in a space \mathbf{U} that is wider than the space $\bar{\mathbf{U}}$ that contains $T\mathbf{u}$. If we have a computationally inexpensive continuous mapping \mathbb{G} such that $\mathbb{G}(T\mathbf{v}_h) \in \bar{\mathbf{U}}, \forall \mathbf{v}_h \in \mathbf{V}_h$. then, probably, the function $\mathbb{G}(T\mathbf{u}_h)$ is much closer to $T\mathbf{u}$ than $T\mathbf{u}_h$.



These arguments form the basis of various **post-processing algorithms** that change a computed solution in accordance with some a priori knowledge of properties of the exact solution. If the error caused by **violations of a priori regularity properties** is dominant and the post-processing operator \mathbb{G} is properly constructed, then

$$\|\mathbb{G}\mathbf{T}\mathbf{u}_h - \mathbf{T}\mathbf{u}\| \ll \|\mathbf{T}\mathbf{u}_h - \mathbf{T}\mathbf{u}\|.$$

In this case, the explicitly computable norm $\|\mathbb{G}\mathbf{T}\mathbf{u}_h - \mathbf{T}\mathbf{u}_h\|$ can be used to evaluate upper and lower bounds of the error.

Indeed, assume that there is a positive number $\alpha < 1$ such that for the mapping \mathbf{T} the estimate

$$\|\mathbb{G}\mathbf{T}\mathbf{u}_h - \mathbf{T}\mathbf{u}\| \leq \alpha \|\mathbf{T}\mathbf{u}_h - \mathbf{T}\mathbf{u}\|.$$

Two-sided estimate

Then, for $\mathbf{e} = \mathbf{u}_h - \mathbf{u}$ we have

$$\begin{aligned}
 (\mathbf{1} - \alpha) \|\mathbf{T}\mathbf{e}\| &= (\mathbf{1} - \alpha) \|\mathbf{T}\mathbf{u}_h - \mathbf{T}\mathbf{u}\| \leq \\
 &\leq \|\mathbf{T}\mathbf{u}_h - \mathbf{T}\mathbf{u}\| - \|\mathbb{G}\mathbf{T}\mathbf{u}_h - \mathbf{T}\mathbf{u}\| \leq \\
 &\leq \|\mathbb{G}\mathbf{T}\mathbf{u}_h - \mathbf{T}\mathbf{u}_h\| \leq \\
 &\leq \|\mathbb{G}\mathbf{T}\mathbf{u}_h - \mathbf{T}\mathbf{u}\| + \|\mathbf{T}\mathbf{u}_h - \mathbf{T}\mathbf{u}\| \leq \\
 &\leq (\mathbf{1} + \alpha) \|\mathbf{T}\mathbf{u}_h - \mathbf{T}\mathbf{u}\| = (\mathbf{1} + \alpha) \|\mathbf{T}\mathbf{e}\|.
 \end{aligned}$$

Thus, if $\alpha \ll 1$, then

$$\|\mathbf{T}\mathbf{u}_h - \mathbf{T}\mathbf{u}\| \simeq \|\mathbb{G}\mathbf{T}\mathbf{u}_h - \mathbf{T}\mathbf{u}_h\|.$$

and the right-hand can be used as an **error indicator**.

Post-processing by averaging

Post-processing operators are often constructed by averaging $T\mathbf{u}_h$ on finite element patches or on the entire domain.

Integral averaging on patches

If $T\mathbf{u}_h \in \mathbf{L}^2$, then post-processing operators are obtained by various averaging procedures. Let Ω_i be a **patch** of M_i elements, i.e.,

$$\overline{\Omega}_i = \bigcup T_{ij}, \quad j = 1, 2, \dots, M_i.$$

Let $\mathbf{P}^k(\Omega_i, \mathbb{R}^n)$ be a subspace of $\overline{\mathbf{U}}$ that consists of vector-valued polynomial functions of degrees less than or equal to \mathbf{k} . Define $\mathbf{g}_i \in \mathbf{P}^k(\Omega_i, \mathbb{R}^n)$ as the minimizer of the problem:

$$\inf_{\mathbf{g} \in \mathbf{P}^k(\Omega_i, \mathbb{R}^n)} \int_{\Omega_i} |\mathbf{g} - T\mathbf{u}_h|^2 \, dx.$$

The minimizer \mathbf{g}_i is used to define the values of an averaged function at some points (nodes). Further, these values are utilized by a prolongation procedure that defines an averaged function

$$\mathbb{G}\mathbf{T}\mathbf{u}_h : \Omega \rightarrow \mathbb{R}.$$

Consider the simplest case. Let \mathbb{T} be the operator ∇ and \mathbf{u}_h be a piecewise affine continuous function. Then,

$$\nabla\mathbf{u}_h \in \mathbf{P}^0(\mathbf{T}_{ij}, \mathbb{R}^n) \quad \text{on each } \mathbf{T}_{ij} \subset \Omega_i.$$

We denote the values of $\nabla\mathbf{u}_h$ on \mathbf{T}_{ij} by $(\nabla\mathbf{u}_h)_{ij}$.

Set $\mathbf{k} = \mathbf{0}$ and find $\mathbf{g}_i \in \mathbf{P}^0$ such that

$$\begin{aligned} \int_{\Omega_i} |\mathbf{g}_i - \nabla \mathbf{u}_h|^2 \, d\mathbf{x} &= \inf_{\mathbf{g} \in \mathbf{P}^0(\Omega_i)} \int_{\Omega_i} |\mathbf{g} - \nabla \mathbf{u}_h|^2 \, d\mathbf{x} = \\ &= \inf_{\mathbf{g} \in \mathbf{P}^0(\Omega_i)} \left\{ |\mathbf{g}|^2 |\Omega_i| - 2\mathbf{g} \cdot \sum_{j=1}^{M_i} (\nabla \mathbf{u}_h)_{ij} |\mathbf{T}_{ij}| + \sum_{j=1}^{M_i} |(\nabla \mathbf{u}_h)_{ij}|^2 |\mathbf{T}_{ij}| \right\}. \end{aligned}$$

It is easy to see that \mathbf{g}_i is given by a weighted sum of $(\nabla \mathbf{u}_h)_{ij}$, namely,

$$\mathbf{g}_i = \sum_{j=1}^{M_i} \frac{|\mathbf{T}_{ij}|}{|\Omega_i|} (\nabla \mathbf{u}_h)_{ij}.$$

Set $\mathbb{G}(\nabla \mathbf{u}_h)(\mathbf{x}_i) = \mathbf{g}_i$.

Repeat this procedure for all nodes and define the vector-valued function $\mathbb{G}\nabla(\mathbf{u}_h)$ by the piecewise affine prolongation of these values. For regular meshes with equal $|\mathbf{T}_{ij}|$, we have

$$\mathbf{g}_i = \sum_{j=1}^{M_i} \frac{1}{M_i} (\nabla \mathbf{u}_h)_{ij}.$$

Various averaging formulas of this type are represented in the form

$$\mathbf{g}_i = \sum_{j=1}^{M_i} \lambda_{ij} (\nabla \mathbf{u}_h)_{ij}, \quad \sum_{j=1}^{M_i} \lambda_{ij} = \mathbf{1},$$

where λ_{ij} are the weight factors. For internal nodes, they may be taken, e.g., as follows

$$\lambda_{ij} = \frac{|\gamma_{ij}|}{2\pi}, \quad |\gamma_{ij}| \text{ is the angle.}$$

However, if a node **belongs to the boundary**, then it is better to choose **special weights**. Their values depend on the mesh and on the type of the boundary. Concerning this point see

I. Hlaváček and M. Krizek. On a superconvergence finite element scheme for elliptic systems. I. Dirichlet boundary conditions. *Aplikace Matematiky*, 32(1987), No.2, 131-154.

Discrete averaging on patches

Consider the problem

$$\inf_{\mathbf{g} \in \mathbb{P}^k(\Omega_i)} \sum_{s=1}^{m_i} |\mathbf{g}(\mathbf{x}_s) - \mathbb{T}\mathbf{u}_h(\mathbf{x}_s)|^2,$$

where the points \mathbf{x}_s are specially selected in Ω_i . Usually, the points \mathbf{x}_s are the so-called **superconvergent points**.

Let $\mathbf{g}_i \in \mathbb{P}^k(\Omega_i)$ be the minimizer of this problem.

If $\mathbf{k} = \mathbf{0}$, and $\mathbb{T} = \nabla$ then

$$\mathbf{g}_i = \frac{1}{m_i} \sum_{s=1}^{m_i} \nabla \mathbf{u}_h(\mathbf{x}_s).$$

Global averaging

Global averaging makes the post-processing not on patches, but on the whole domain.

Assume that $\mathbf{T}\mathbf{u}_h \in \mathbf{L}^2$ and find $\bar{\mathbf{g}}_h \in \mathbf{V}_h(\Omega) \subset \bar{\mathbf{U}}$ such that

$$\|\bar{\mathbf{g}}_h - \mathbf{T}\mathbf{u}_h\|_{\Omega}^2 = \inf_{\mathbf{g}_h \in \mathbf{V}_h(\Omega)} \|\mathbf{g}_h - \mathbf{T}\mathbf{u}_h\|_{\Omega}^2.$$

The function $\bar{\mathbf{g}}_h$ can be viewed as $\mathbb{G}\mathbf{T}\mathbf{u}_h$. Very often $\bar{\mathbf{g}}_h$ is a better image of $\mathbf{T}\mathbf{u}$ than the functions obtained by local procedures.

Remark

Moreover, mathematical justifications of the methods based on global averaging procedures can be performed under weaker assumptions what makes them applicable to a wider class of problems see, e.g.,

Carstensen, C.; Bartels, S. Each averaging technique yields reliable a posteriori error control in FEM on unstructured grids. I: Low order conforming, nonconforming, and mixed FEM, *Math. Comp.*, 71(2002)

Task 5

Solve the boundary-value problem

$$\Delta \mathbf{u} + \mathbf{f} = \mathbf{0}, \quad \mathbf{u} = \mathbf{0} \text{ on } \partial\Omega$$

by h -version FEM (use Matlab or another code). Apply the simplest gradient-averaging error indicator to indicate the error distribution. Compare it with the distribution of true error (the latter can be extracted from a solution on a much finer mesh).

Justifications of the method. Superconvergence

Let \mathbf{u}_h be a Galerkin approximation of \mathbf{u} computed on \mathbf{V}_h . For piecewise affine approximations of the diffusion problem we have the estimate

$$\|\nabla(\mathbf{u} - \mathbf{u}_h)\|_{2,\Omega} \leq \mathbf{c}_1 \mathbf{h}, \quad \|\mathbf{u} - \mathbf{u}_h\|_{2,\Omega} \leq \mathbf{c}_2 \mathbf{h}^2$$

However, it was discovered see, e.g.,

L. A. Oganjesjan and L. A. Ruchovec. *Z. Vjycisl. Mat. i Mat. Fiz.*,9(1969);

M. Zlámal. *Lecture Notes*. Springer, 1977;

L. B. Wahlbin. *Lecture Notes*. Springer, 1969 that in certain cases **this rate may be higher**. For example it may happen that

$$|\mathbf{u}(\mathbf{x}_s) - \mathbf{u}_h(\mathbf{x}_s)| \leq \mathbf{C} \mathbf{h}^{2+\sigma} \quad \sigma > 0$$

at a **superconvergent point** \mathbf{x}_s .

Certainly, existence and location of superconvergent points strongly depends on the structure of \mathcal{T}_h .

For the paradigm of the diffusion problem we say that an operator \mathbb{G} possesses a *superconvergence* property in $\omega \subset \Omega$ if

$$\|\nabla \mathbf{u} - \mathbb{G} \nabla \mathbf{u}_h\|_{2,\omega} \leq \mathbf{c}_2 \mathbf{h}^{1+\sigma},$$

where the constant \mathbf{c}_2 may depend on higher norms of \mathbf{u} and the structure of \mathcal{T}_h .

For the diffusion problem estimates of such a type can be found, e.g., in
I. Hlaváček and M. Krizek. On a superconvergence finite element scheme
for elliptic systems. I. Dirichlet boundary conditions. *Aplikace Matematiky*,
32(1987).

M. Krížek and P. Neittaanmäki. Superconvergence phenomenon in the
finite element method arising from averaging of gradients *Numer. Math.*,
45(1984)

By exploiting the superconvergence properties, e.g.,

$$\|\nabla \mathbf{u} - \mathbb{G}\nabla \mathbf{u}_h\|_{2,\omega} \leq \mathbf{c}_2 \mathbf{h}^{1+\sigma},$$

while

$$\|\nabla \mathbf{u} - \nabla \mathbf{u}_h\|_{2,\omega} \leq \mathbf{c}_2 \mathbf{h},$$

one can usually construct a simple post-processing operator \mathbb{G} satisfying the condition

$$\|\mathbb{G}\nabla \mathbf{u}_h - \nabla \mathbf{u}\| \leq \alpha \|\nabla \mathbf{u}_h - \nabla \mathbf{u}\|.$$

where the value of α decreases as \mathbf{h} tends to zero.

Since

$$\begin{aligned}\|\mathbb{G}\nabla\mathbf{u}_h - \nabla\mathbf{u}_h\| &\leq \|\nabla\mathbf{u}_h - \nabla\mathbf{u}\| + \|\mathbb{G}\nabla\mathbf{u}_h - \nabla\mathbf{u}\|, \\ \|\mathbb{G}\nabla\mathbf{u}_h - \nabla\mathbf{u}_h\| &\geq \|\nabla\mathbf{u}_h - \nabla\mathbf{u}\| - \|\mathbb{G}\nabla\mathbf{u}_h - \nabla\mathbf{u}\|.\end{aligned}$$

where the first term in the right-hand side is of the order \mathbf{h} and the second one is of $\mathbf{h}^{1+\delta}$. We see that

$$\|\mathbb{G}\nabla\mathbf{u}_h - \nabla\mathbf{u}_h\| \sim \mathbf{h}$$

Therefore, we observe that in the decomposition

$$\|\nabla(\mathbf{u}_h - \mathbf{u})\| \leq \|\nabla\mathbf{u}_h - \mathbb{G}\nabla\mathbf{u}_h\| + \|\mathbb{G}\nabla\mathbf{u}_h - \nabla\mathbf{u}\|$$

asymptotically dominates the second directly computable term.

Thus, we obtain a simple error indicator:

$$\|\nabla(\mathbf{u}_h - \mathbf{u})\| \approx \|\nabla\mathbf{u}_h - \mathbb{G}\nabla\mathbf{u}_h\|.$$

Note that

$$\mathbf{i}_{\text{eff}} = \frac{\|\nabla(\mathbf{u}_h - \mathbf{u})\|}{\|\nabla\mathbf{u}_h - \mathbb{G}\nabla\mathbf{u}_h\|} \approx \mathbf{1} + \mathbf{c}\mathbf{h}^\delta$$

so that this error indicator is **asymptotically exact** provided that \mathbf{u}_h is a Galerkin approximation, \mathbf{u} is sufficiently regular and \mathbf{h} is small enough. Such type error indicators (often called **ZZ indicators** by the names of Zienkiewicz and Zhu) are widely used as cheap error indicators in engineering computations.

Some references

M. Ainsworth, J. Z. Zhu, A. W. Craig and O. C. Zienkiewicz. Analysis of the Zienkiewicz-Zhu a posteriori error estimator in the finite element method, *Int. J. Numer. Methods Engrg.*, 28(1989).

I. Babuska and R. Rodriguez. The problem of the selection of an a posteriori error indicator based on smoothing techniques, *Internat. J. Numer. Meth. Engrg.*, 36(1993).

O. C. Zienkiewicz and J. Z. Zhu. A simple error estimator and adaptive procedure for practical engineering analysis, *Internat. J. Numer. Meth. Engrg.*, 24(1987)

Post-processing by equilibration

For a solution of the diffusion problem we know that

$$\mathbf{div}\boldsymbol{\sigma} + \mathbf{f} = \mathbf{0},$$

where $\boldsymbol{\sigma} = \mathbf{A}\nabla\mathbf{u}$. This suggests an idea to construct an operator \mathbb{G} such that

$$\mathbf{div}(\mathbb{G}(\mathbf{A}\nabla\mathbf{u}_h)) + \mathbf{f} = \mathbf{0}.$$

If \mathbb{G} possesses additional properties (linearity, boundedness), then we may hope that the function $\mathbb{G}\mathbf{A}\nabla\mathbf{u}_h$ is closer to $\boldsymbol{\sigma}$ than $\mathbf{A}\nabla\mathbf{u}_h$ and use the quantity $\|\mathbf{A}\nabla\mathbf{u}_h - \mathbb{G}\mathbf{A}\nabla\mathbf{u}_h\|$ as an error indicator.

This idea can be applied to an important class of problems

$$\mathbf{\Lambda}^* \mathbb{T}u + f = 0, \quad \mathbb{T}u = \mathcal{A}\mathbf{\Lambda}u, \quad (56)$$

where \mathcal{A} is a positive definite operator, $\mathbf{\Lambda}$ is a linear continuous operator, and $\mathbf{\Lambda}^*$ is the adjoint operator.

In continuum mechanics, equations of the type (56) are referred to as the **equilibrium equations**. Therefore, it is natural to call an operator \mathbb{G} an **equilibration** operator.

If the equilibration has been performed exactly then it is not difficult to get an upper error bound. However, in general, this task is either cannot be fulfilled or lead to complicated and expensive procedures. Known methods are usually end with approximately equilibrated fluxes.

Goal-oriented error estimates

Global error estimates give a general idea on the quality of an approximate solution and stopping criteria. However, often it is useful to estimate the errors in terms of **specially selected linear functionals** ℓ_s , $s = 1, 2, \dots, M$, e.g.,

$$\langle \ell, \mathbf{v} - \mathbf{u} \rangle = \int_{\Omega} \phi_0 (\mathbf{v} - \mathbf{u}) \, d\mathbf{x},$$

where ϕ is a locally supported function. Since

$$| \langle \ell, \mathbf{u} - \mathbf{u}_h \rangle | \leq \| \ell \| \| \mathbf{u} - \mathbf{u}_h \|_{\mathbf{v}},$$

we can obtain such an estimate throughout the global a posteriori estimate. However, in many cases, such a method will strongly overestimate the quantity.

Adjoint problem

A posteriori estimates of the errors evaluated in terms of linear functionals are derived by attracting the **adjoint** boundary-value problem whose right-hand side is formed by the functional ℓ .

Let us represent this idea in the simplest form. Consider a system

$$\mathbf{A}\mathbf{u} = \mathbf{f},$$

where \mathbf{A} is a positive definite matrix and \mathbf{f} is a given vector. Let \mathbf{v} be an approximate solution. Define \mathbf{u}_ℓ by the relation

$$\mathbf{A}^*\mathbf{u}_\ell = \ell,$$

where \mathbf{A}^* is the matrix adjoint to \mathbf{A} . Then,

$$\ell \cdot (\mathbf{u} - \mathbf{v}) = \mathbf{A}^*\mathbf{u}_\ell \cdot \mathbf{u} - \ell \cdot \mathbf{v} = \mathbf{f} \cdot \mathbf{u}_\ell - \ell \cdot \mathbf{v} = (\mathbf{f} - \mathbf{A}\mathbf{v}) \cdot \mathbf{u}_\ell$$

Certainly, the above consideration holds in a more general (operator) sense, so that for a pair of operators A and A^* we have

$$\langle \ell, \mathbf{u} - \mathbf{v} \rangle = \langle \mathbf{f} - A\mathbf{v}, \mathbf{u}_\ell \rangle. \quad (57)$$

and find the error with respect to a linear functional by the product of the **residual** and the **exact solution of the adjoint problem**:

$$A^* \mathbf{u}_\ell = \ell.$$

Practical application of this principle depends on the ability to find either \mathbf{u}_ℓ or its sharp approximation.

Consider again the diffusion problem. Now, it is convenient to denote the solution of the original problem by \mathbf{u}_f , i.e

$$\int_{\Omega} \mathbf{A} \nabla \mathbf{u}_f \cdot \nabla \mathbf{w} \, d\mathbf{x} = \int_{\Omega} \mathbf{f} \mathbf{w} \, d\mathbf{x}, \quad \forall \mathbf{w} \in \mathbf{V}_0(\Omega).$$

Since in our case $\mathbf{A} = \mathbf{A}^*$, the **adjoint** problem is to find $\mathbf{u}_\ell \in \mathbf{V}_0(\Omega)$ such that

$$\int_{\Omega} \mathbf{A} \nabla \mathbf{u}_\ell \cdot \nabla \mathbf{w} \, d\mathbf{x} = \int_{\Omega} \ell \mathbf{w} \, d\mathbf{x}, \quad \forall \mathbf{w} \in \mathbf{V}_0(\Omega).$$

Let Ω be divided into a number of elements \mathbf{T}_i , $i = 1, 2, \dots, N$. Given approximations on the elements, we define a finite-dimensional subspace $\mathbf{V}_{0h} \in \mathbf{V}_0(\Omega)$ and the Galerkin approximations \mathbf{u}_{fh} and $\mathbf{u}_{\ell h}$:

$$\int_{\Omega} \mathbf{A} \nabla \mathbf{u}_{fh} \cdot \nabla \mathbf{w}_h \, d\mathbf{x} = \int_{\Omega} \mathbf{f} \mathbf{w}_h \, d\mathbf{x}, \quad \forall \mathbf{w}_h \in \mathbf{V}_{0h},$$

$$\int_{\Omega} \mathbf{A} \nabla \mathbf{u}_{\ell h} \cdot \nabla \mathbf{w}_h \, d\mathbf{x} = \int_{\Omega} \ell \mathbf{w}_h \, d\mathbf{x}, \quad \forall \mathbf{w}_h \in \mathbf{V}_{0h}.$$

Since

$$\int_{\Omega} \ell(\mathbf{u}_f - \mathbf{u}_{fh}) \, d\mathbf{x} = \int_{\Omega} \mathbf{A} \nabla \mathbf{u}_{\ell} \cdot \nabla(\mathbf{u}_f - \mathbf{u}_{fh}) \, d\mathbf{x}$$

and

$$\int_{\Omega} \mathbf{A} \nabla \mathbf{u}_{\ell h} \cdot \nabla(\mathbf{u}_f - \mathbf{u}_{fh}) \, d\mathbf{x} = 0,$$

We arrive at the relation

$$\int_{\Omega} \ell(\mathbf{u}_f - \mathbf{u}_{fh}) \, d\mathbf{x} = \int_{\Omega} \mathbf{A} \nabla(\mathbf{u}_\ell - \mathbf{u}_{\ell h}) \cdot \nabla(\mathbf{u}_f - \mathbf{u}_{fh}) \, d\mathbf{x} \quad (58)$$

whose right-hand side is expressed in the form

$$\begin{aligned} \sum_{i=1}^N \int_{T_i} \mathbf{A} \nabla(\mathbf{u}_f - \mathbf{u}_{fh}) \cdot \nabla(\mathbf{u}_\ell - \mathbf{u}_{\ell h}) \, d\mathbf{x} = \\ \sum_{i=1}^N \left\{ - \int_{T_i} \operatorname{div}(\mathbf{A} \nabla(\mathbf{u}_f - \mathbf{u}_{fh})) (\mathbf{u}_\ell - \mathbf{u}_{\ell h}) \, d\mathbf{x} + \right. \\ \left. + \frac{1}{2} \int_{\partial T_i} \mathbf{j}(\nu_i \cdot \mathbf{A} \nabla(\mathbf{u}_f - \mathbf{u}_{fh})) (\mathbf{u}_\ell - \mathbf{u}_{\ell h}) \, d\mathbf{s} \right\}. \end{aligned}$$

This relation implies the estimate

$$\begin{aligned}
\int_{\Omega} \ell(\mathbf{u}_f - \mathbf{u}_{fh}) \mathbf{d}\mathbf{x} &= \sum_{i=1}^N \left\{ \|\mathbf{div} \mathbf{A} \nabla(\mathbf{u}_f - \mathbf{u}_{fh})\|_{2, \mathcal{T}_i} \|\mathbf{u}_\ell - \mathbf{u}_{\ell h}\|_{2, \mathcal{T}_i} + \right. \\
&\quad \left. + \frac{1}{2} \|\mathbf{j}(\boldsymbol{\nu}_i \cdot \mathbf{A} \nabla(\mathbf{u}_f - \mathbf{u}_{fh}))\|_{2, \partial \mathcal{T}_i} \|\mathbf{u}_\ell - \mathbf{u}_{\ell h}\|_{2, \partial \mathcal{T}_i} \right\} = \\
&= \sum_{i=1}^N \left\{ \|\mathbf{f} + \mathbf{div} \mathbf{A} \nabla \mathbf{u}_{fh}\|_{2, \mathcal{T}_i} \|\mathbf{u}_\ell - \mathbf{u}_{\ell h}\|_{2, \mathcal{T}_i} + \right. \\
&\quad \left. + \frac{1}{2} \|\mathbf{j}(\boldsymbol{\nu}_i \cdot \mathbf{A} \nabla \mathbf{u}_{fh})\|_{2, \partial \mathcal{T}_i} \|\mathbf{u}_\ell - \mathbf{u}_{\ell h}\|_{2, \partial \mathcal{T}_i} \right\}.
\end{aligned}$$

Here, the principal terms are the same as in the explicit residual method, but the weights are given by the norms of $\mathbf{u}_\ell - \mathbf{u}_{\ell h}$.

Assume that $\mathbf{u}_\ell \in \mathbf{H}^2$ and $\mathbf{u}_{\ell h}$ is constructed by piecewise affine continuous approximations. Then the norms $\|\mathbf{u}_\ell - \mathbf{u}_{\ell h}\|_{\mathbf{T}_i}$ and $\|\mathbf{u}_\ell - \mathbf{u}_{\ell h}\|_{2, \partial \mathbf{T}_i}$ are estimated by the quantities $h^\alpha |\mathbf{u}_\ell|_{2,2, \mathbf{T}_i}$ with $\alpha = 1$ and $1/2$ and the multipliers $\hat{\mathbf{c}}_i$ and $\hat{\mathbf{c}}_{ij}$, respectively.

In this case, we obtain an estimate with constants defined by the standard

$$\mathbf{H}^2 \rightarrow \mathbf{V}_{0h}$$

interpolation operator whose evaluation is much simpler than that of the constants arising in the

$$\mathbf{H}^1 \rightarrow \mathbf{V}_{0h}$$

interpolation.

A posteriori estimates in L^2 -norm

In principle, this technology can be exploited to evaluate estimates in L^2 -norm. Indeed,

$$\begin{aligned}
 \|\mathbf{u}_f - \mathbf{u}_{fh}\| &= \sup_{\boldsymbol{\ell} \in L^2} \frac{(\boldsymbol{\ell}, \mathbf{u}_f - \mathbf{u}_{fh})}{\|\boldsymbol{\ell}\|} = \sup_{\boldsymbol{\ell} \in L^2} \frac{(\mathbf{A}\nabla\mathbf{u}_\ell, \nabla(\mathbf{u}_f - \mathbf{u}_{fh}))}{\|\boldsymbol{\ell}\|} = \\
 &= \sup_{\boldsymbol{\ell} \in L^2} \frac{(\mathbf{A}\nabla(\mathbf{u}_\ell - \pi_h(\mathbf{u}_\ell)), \nabla(\mathbf{u}_f - \mathbf{u}_{fh}))}{\|\boldsymbol{\ell}\|} = \\
 &= \sup_{\boldsymbol{\ell} \in L^2} \frac{(\nabla(\mathbf{u}_\ell - \pi_h(\mathbf{u}_\ell)), \mathbf{A}\nabla(\mathbf{u}_f - \mathbf{u}_{fh}))}{\|\boldsymbol{\ell}\|} = \\
 &= \sup_{\boldsymbol{\ell} \in L^2} \frac{\sum_{i=1}^N \left\{ \int_{T_i} \nabla(\mathbf{u}_\ell - \pi_h(\mathbf{u}_\ell)), \mathbf{A}\nabla(\mathbf{u}_f - \mathbf{u}_{fh}) \, dx \right\}}{\|\boldsymbol{\ell}\|}
 \end{aligned}$$

Integrating by parts, we obtain

$$\frac{\sum_{i=1}^N \left\{ \|\mathbf{f} + \mathbf{div} \mathbf{A} \nabla \mathbf{u}_{\text{fh}}\|_{\mathbf{T}_i} \|\mathbf{u}_\ell - \pi_{\text{h}}(\mathbf{u}_\ell)\|_{\mathbf{T}_i} + \frac{1}{2} \|\mathbf{j}(\boldsymbol{\nu}_i \cdot \mathbf{A} \nabla \mathbf{u}_{\text{fh}})\|_{\partial \mathbf{T}_i} \|\mathbf{u}_\ell - \pi_{\text{h}}(\mathbf{u}_\ell)\|_{\partial \mathbf{T}_i} \right\}}{\|\boldsymbol{\ell}\|}$$

If for *any* $\boldsymbol{\ell} \in L^2$ the adjoint problem has a regular solution (e.g., $\mathbf{u}_\ell \in \mathbf{H}^2$), so that we could combine the standard interpolation estimate for the interpolant of \mathbf{u}_ℓ with the regularity estimate for the PDE (e.g., $\|\mathbf{u}_\ell\| \leq \mathbf{C}_1 \|\boldsymbol{\ell}\|$), then we obtain

$$\|\mathbf{u}_\ell - \pi_{\text{h}}(\mathbf{u}_\ell)\|_{\mathbf{T}_i} \leq \mathbf{C}_1 \mathbf{h}^{\alpha_1} \|\boldsymbol{\ell}\|, \quad \|\mathbf{u}_\ell - \pi_{\text{h}}(\mathbf{u}_\ell)\|_{\partial \mathbf{T}_i} \leq \mathbf{C}_1 \mathbf{h}^{\alpha_2} \|\boldsymbol{\ell}\|$$

with certain $\alpha_{\mathbf{k}}$.

Under the above conditions $\|\boldsymbol{\ell}\|$ is reduced and we arrive at the estimate, in which the element residuals and interelement jumps are weighted with factors $\mathbf{C}_1 \mathbf{h}^{\alpha_1}$ and $\mathbf{C}_2 \mathbf{h}^{\alpha_2}$.

References

Methods using adjoint problems has been investigated in the works of R. Becker, C. Johnson, R. Rannacher and other scientists. A more detailed exposition of these works can be found in

W. Bangerth and R. Rannacher. *Adaptive finite element methods for differential equations*. Birkhäuser, Berlin, 2003.

R. Becker and R. Rannacher. A feed-back approach to error control in finite element methods: Basic approach and examples, *East-West J. Numer. Math.*, 4(1996), 237-264.

Concerning error estimation in goal-oriented quantities we refer, e.g., to

J. T. Oden, S. Prudhomme. Goal-oriented error estimation and adaptivity for the finite element method, *Comput. Math. Appl.*, 41, 735-756, 2001.

S. Korotov, P. Neittaanmaki and S. Repin. A posteriori error estimation of goal-oriented quantities by the superconvergence patch recovery, *J. Numer. Math.* 11 (2003)

Comment

In the literature devoted to a posteriori error analysis one can find often find terms like

"duality approach to a posteriori error estimation" or
"dual-based error estimates".

However, the essence behind such a terminology may be quite different because the word *"duality"* is used in 3 different meanings:
(a) Duality in the sense of functional spaces. We have seen that if for the equation $\mathcal{L}\mathbf{u} = \mathbf{f}$ errors are measured in the original (energy) norm then a consistent upper bound is given by the residual in the norm of the space **topologically dual** to a subspace of the energy space (e.g., \mathbf{H}^{-1}).

(b) Duality in the sense of using the Adjoint Problem.

(c) Duality in the sense of the Theory of the Calculus of Variations.

**In the next lecture
we will proceed to the detailed exposition
of the approach (c).**

Lecture 3

In the lecture, we derive Functional A Posteriori Estimate for the problem

$$\Delta \mathbf{u} + \mathbf{f} = \mathbf{0}, \Omega \quad \mathbf{u} = \mathbf{0} \partial\Omega.$$

and discuss its meaning, principal features and practical implementation.

Lecture plan

1. **Functional a posteriori estimates.**
2. How to derive them? Paradigm of a simple elliptic problem
3. How to use them in practice?
4. Examples.

Lecture plan

1. Functional a posteriori estimates.
2. **How to derive them? Paradigm of a simple elliptic problem**
3. How to use them in practice?
4. Examples.

Lecture plan

1. Functional a posteriori estimates.
2. How to derive them? Paradigm of a simple elliptic problem
3. **How to use them in practice?**
4. Examples.

Lecture plan

1. Functional a posteriori estimates.
2. How to derive them? Paradigm of a simple elliptic problem
3. How to use them in practice?
4. **Examples.**

Functional A Posteriori Estimates

Functional A Posteriori Estimate is a computable majorant of the difference between exact solution \mathbf{u} and any conforming approximation \mathbf{v} having the general form:

$$\Phi(\mathbf{u} - \mathbf{v}) \leq M(\mathcal{D}, \mathbf{v}) \quad \forall \mathbf{v} \in \mathbf{V}! \quad (59)$$

\mathcal{D} is the data set (coefficients, domain, parameters, etc.),
 $\Phi : \mathbf{V} \rightarrow \mathbb{R}_+$ is a given functional.
 M must be computable and continuous in the sense that

$$M(\mathcal{D}, \mathbf{v}) \rightarrow 0, \quad \text{if } \mathbf{v} \rightarrow \mathbf{u}$$

Types of Φ

- Energy norm $\Phi(\mathbf{u} - \mathbf{v}) = \|\mathbf{u} - \mathbf{v}\|_{\Omega}$
- Local norm $\Phi(\mathbf{u} - \mathbf{v}) = \|\mathbf{u} - \mathbf{v}\|_{\omega}$
- Goal-oriented quantity $\Phi(\mathbf{u} - \mathbf{v}) = (\ell, \mathbf{u} - \mathbf{v})$

Functional a posteriori estimate gives complete solution of the error control problem from the viewpoint of the **MATHEMATICAL THEORY** of PDE's

METHODS OF THE DERIVATION.

These estimates are derived by purely functional methods using the analysis of variational problems or integral identities.

Variational method 96'-97'

Exploits variational structure of the original problem and Duality Theory in the Calculus of Variations.

See

S. Repin *Mathematics of Computation*, 69(230), pp. 2000, 481-500.

A systematic exposition of the variational approach to deriving Functional a Posteriori Estimates can be found in

P. Neittaanmaki and S. Repin. *Reliable methods for computer simulation. Error control and a posteriori estimates*. Elsevier, NY, 2004

Nonvariational method 2000'

Derives a posteriori estimates by certain transformations of integral identities.

Basic idea of the method is presented in

S. Repin. *Proc. St.-Petersburg Math. Society*, 2001 pp. 148-179 (in Russian, translated in *American Mathematical Translations Series 2*, 9(2003)). .

Other publications:

S. Repin. Estimates of deviation from exact solutions of initial-boundary value problems for the heat equation, *Rend. Mat. Acc. Lincei*, 13(2002), pp. 121-133.

S. Repin *Estimates of deviations from exact solutions for some boundary-value problems with incompressibility condition*, *Algebra and Analiz*, 16(2004), 5, pp. 124-161.

Let us consider both methods in application to our basic problem

Variational Method

Let \mathbf{u} be a (generalized) solution of the problem

$$\Delta \mathbf{u} + \mathbf{f} = \mathbf{0}, \quad \Omega \quad \mathbf{u} = \mathbf{0} \quad \partial\Omega.$$

As we have seen in Lecture 1, this problem is equivalent to the following variational problem:

Problem \mathcal{P} . Find $\mathbf{u} \in \mathbf{V}_0 := \mathring{\mathbf{H}}^1(\Omega)$ such that

$$\mathbf{J}(\mathbf{u}) = \inf_{\mathbf{v} \in \mathbf{V}_0} \mathbf{J}(\mathbf{v}),$$

where

$$\mathbf{J}(\mathbf{v}) = \frac{1}{2} \|\nabla \mathbf{v}\|^2 - (\mathbf{f}, \mathbf{v}).$$

By the reasons that we discussed earlier this problem has a unique solution.

Lagrangian

Note that

$$J(\mathbf{v}) = \sup_{\mathbf{y} \in \mathbf{Y}} L(\nabla \mathbf{v}, \mathbf{y}), \quad L(\nabla \mathbf{v}, \mathbf{y}) = \int_{\Omega} \left(\nabla \mathbf{v} \cdot \mathbf{y} - \frac{1}{2} |\mathbf{y}|^2 - f\mathbf{v} \right) dx$$

where $\mathbf{Y} = \mathbf{L}^2(\Omega, \mathbb{R}^n)$. Indeed, the value of the above supremum cannot exceed the one we obtain if for almost all $\mathbf{x} \in \Omega$ solve the pointwise problems

$$\sup_{\mathbf{y}(\mathbf{x})} (\nabla \mathbf{v})(\mathbf{x}) \cdot \mathbf{y}(\mathbf{x}) - \frac{1}{2} |\mathbf{y}(\mathbf{x})|^2 \quad \mathbf{x} \in \Omega$$

whose upper bound is attained if set $\mathbf{y}(\mathbf{x}) = (\nabla \mathbf{v})(\mathbf{x})$. Since $\nabla \mathbf{v} \in \mathbf{Y}$, we observe that the respective maximizer belongs to \mathbf{Y} and, therefore

$$\sup_{\mathbf{y} \in \mathbf{Y}} L(\nabla \mathbf{v}, \mathbf{y}) = L(\nabla \mathbf{v}, \nabla \mathbf{v}) = J(\mathbf{v}).$$

Minimax Formulations

Then, the original problem comes in the **minimax** form:

$$(\mathcal{P}) \quad \inf_{\mathbf{v} \in \mathbf{V}_0} \sup_{\mathbf{y} \in \mathbf{Y}} \mathbf{L}(\nabla \mathbf{v}, \mathbf{y})$$

If the order of **inf** and **sup** is changed, then we arrive at the so-called **dual problem**

$$(\mathcal{P}^*) \quad \sup_{\mathbf{y} \in \mathbf{Y}} \inf_{\mathbf{v} \in \mathbf{V}_0} \mathbf{L}(\nabla \mathbf{v}, \mathbf{y})$$

Note that

$$\begin{aligned} \inf_{\mathbf{v} \in \mathbf{V}_0} \int_{\Omega} \left(\nabla \mathbf{v} \cdot \mathbf{y} - \frac{1}{2} |\mathbf{y}|^2 - \mathbf{f} \mathbf{v} \right) \mathbf{d}\mathbf{x} &= -\frac{1}{2} \|\mathbf{y}\|^2 + \inf_{\mathbf{v} \in \mathbf{V}_0} \int_{\Omega} (\nabla \mathbf{v} \cdot \mathbf{y} - \mathbf{f} \mathbf{v}) \mathbf{d}\mathbf{x} = \\ &= \begin{cases} -\frac{1}{2} \|\mathbf{y}\|^2 & \text{if } \mathbf{y} \in \mathbf{Q}_f := \{\mathbf{y} \in \mathbf{Y} \mid \operatorname{div} \mathbf{y} + \mathbf{f} = 0\} \\ -\infty & \text{if } \mathbf{y} \notin \mathbf{Q}_f \end{cases} \end{aligned}$$

Dual Problem

Thus, we observe that the dual problem has the form: find $\mathbf{p} \in \mathbf{Q}_f$ such that

$$-I^*(\mathbf{p}) = \sup_{\mathbf{y} \in \mathbf{Q}_f} -I^*(\mathbf{y})$$

where

$$I^*(\mathbf{q}) = \frac{1}{2} \|\mathbf{q}\|^2$$

How are these two problems related?

First, we establish one relation that holds regardless of the structure of the Lagrangian.

Sup Inf and Inf Sup

Lemma

Let $L(x, y)$ be a functional defined on the elements of two nonempty sets X and Y . Then

$$\sup_{y \in Y} \inf_{x \in X} L(x, y) \leq \inf_{x \in X} \sup_{y \in Y} L(x, y). \quad (60)$$

Proof

It is easy to see that

$$L(x, y) \geq \inf_{\xi \in X} L(\xi, y), \quad \forall x \in X, y \in Y.$$

Taking the supremum over $y \in Y$, we obtain

proof

$$\sup_{y \in Y} L(\mathbf{x}, y) \geq \sup_{y \in Y} \inf_{\xi \in X} L(\xi, y), \quad \forall \mathbf{x} \in \mathbf{X}.$$

The left-hand side depends on \mathbf{x} , while the right-hand side is a number. Thus, we may take infimum over $\mathbf{x} \in \mathbf{X}$ and obtain the inequality

$$\inf_{\mathbf{x} \in \mathbf{X}} \sup_{y \in Y} L(\mathbf{x}, y) \geq \sup_{y \in Y} \inf_{\xi \in X} L(\xi, y).$$

Therefore, we always have

$$\sup \mathcal{P}^* \leq \inf \mathcal{P}$$

Duality relations

However, in our case we have a stronger relation, namely

$$\sup \mathcal{P}^* = \inf \mathcal{P}$$

To prove this fact, we note that

$$\int_{\Omega} \nabla \mathbf{u} \cdot \nabla \mathbf{v} \, d\mathbf{x} = \int_{\Omega} \mathbf{f} \mathbf{v} \, d\mathbf{x} \quad \forall \mathbf{v} \in \mathbf{V}_0.$$

Therefore $\mathbf{p} = \nabla \mathbf{u} \in \mathbf{Q}_f$ and

$$-I^*(\mathbf{p}) = -\frac{1}{2} \|\nabla \mathbf{u}\|^2 = \int_{\Omega} \left(\frac{1}{2} |\nabla \mathbf{u}|^2 - |\nabla \mathbf{u}|^2 \right) d\mathbf{x} = \int_{\Omega} \left(\frac{1}{2} |\nabla \mathbf{u}|^2 - \mathbf{f} \mathbf{u} \right) d\mathbf{x} = \mathbf{J}(\mathbf{u}).$$

Let us use the Mikhlin's estimate established in Lecture 2:

$$\frac{1}{2} \|\nabla(\mathbf{u} - \mathbf{v})\|^2 \leq \mathbf{J}(\mathbf{v}) - \mathbf{J}(\mathbf{u}).$$

Since $\mathbf{J}(\mathbf{u}) = -\mathbf{I}^*(\mathbf{p})$, we have

$$\frac{1}{2} \|\nabla(\mathbf{u} - \mathbf{v})\|^2 \leq \mathbf{J}(\mathbf{v}) + \mathbf{I}^*(\mathbf{p}) \leq \mathbf{J}(\mathbf{v}) + \mathbf{I}^*(\mathbf{q}) \quad \forall \mathbf{q} \in \mathbf{Q}_f.$$

Reform this estimate by using the fact that $\mathbf{q} \in \mathbf{Q}_f$.

$$\begin{aligned} \mathbf{J}(\mathbf{v}) + \mathbf{I}^*(\mathbf{q}) &= \frac{1}{2} \|\nabla \mathbf{v}\|^2 - (\mathbf{f}, \mathbf{v}) + \frac{1}{2} \|\mathbf{q}\|^2 \\ &= \frac{1}{2} \|\nabla \mathbf{v}\|^2 + \frac{1}{2} \|\mathbf{q}\|^2 - (\nabla \mathbf{v}, \mathbf{q}) = \\ &= \frac{1}{2} \|\nabla \mathbf{v} - \mathbf{q}\|^2. \end{aligned}$$

Now, we have

$$\|\nabla(\mathbf{v} - \mathbf{u})\| \leq \|\nabla\mathbf{v} - \mathbf{q}\| \quad \forall \mathbf{q} \in \mathbf{Q}_f.$$

Take arbitrary $\mathbf{y} \in \mathbf{L}^2(\Omega)$. Then,

$$\|\nabla(\mathbf{v} - \mathbf{u})\| \leq \|\nabla\mathbf{v} - \mathbf{y}\| + \inf_{\mathbf{q} \in \mathbf{Q}_f} \|\mathbf{y} - \mathbf{q}\|.$$

How to estimate the above infimum?

Various methods give one and the same answer:

$$\inf_{\mathbf{q} \in \mathbf{Q}_f} \|\mathbf{y} - \mathbf{q}\| \leq \mathbf{I} \mathbf{div} + \mathbf{f} \mathbf{I} \quad \mathbf{y} \in \mathbf{L}^2(\Omega), \quad (61)$$

$$\inf_{\mathbf{q} \in \mathbf{Q}_f} \|\mathbf{y} - \mathbf{q}\| \leq \mathbf{C}_\Omega \|\mathbf{div} + \mathbf{f}\| \quad \mathbf{y} \in \mathbf{H}(\Omega, \mathbf{div}), \quad (62)$$

Proof

To prove these estimates we consider an auxiliary problem

$$\Delta \eta + \mathbf{f} + \operatorname{div} \mathbf{y} = 0 \quad \Omega \quad \eta = 0 \quad \partial \Omega.$$

$$\int_{\Omega} \nabla \eta \cdot \nabla \mathbf{w} \, dx = \int_{\Omega} (\mathbf{f} \mathbf{w} - \mathbf{y} \cdot \nabla \mathbf{w}) \, dx$$

\bar{q}

$$\int_{\Omega} \overbrace{(\nabla \eta + \mathbf{y})} \cdot \nabla \mathbf{w} \, dx = \int_{\Omega} \mathbf{f} \mathbf{w} \, dx \quad \forall \mathbf{w} \in \mathbf{V}_0$$

Thus, $\bar{q} \in Q_f$!!!

Since $\boldsymbol{\eta}$ is a solution of the boundary-value problem with right-hand side $\mathbf{div} \mathbf{y} + \mathbf{f} \in \mathbf{H}^{-1}$, we have

$$\|\nabla \boldsymbol{\eta}\| \leq \mathbf{I} \mathbf{div} \mathbf{y} + \mathbf{f} \mathbf{I},$$

Then

$$\inf_{\mathbf{q} \in \mathbf{Q}_f} \|\mathbf{y} - \mathbf{q}\| \leq \|\mathbf{y} - \bar{\mathbf{q}}\| = \|\nabla \boldsymbol{\eta}\| \leq \mathbf{I} \mathbf{div} \mathbf{y} + \mathbf{f} \mathbf{I}.$$

Here

$$\mathbf{I} \mathbf{div} \mathbf{y} + \mathbf{f} \mathbf{I} = \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\int_{\Omega} (\mathbf{y} \cdot \nabla \mathbf{w} - \mathbf{f} \mathbf{w}) \, dx}{\|\nabla \mathbf{w}\|}$$

$\mathbf{y} \in \mathbf{H}(\Omega, \text{div})$

If \mathbf{y} has a square summable divergence, then we have

$$\|\mathbf{div} \mathbf{y} + \mathbf{f}\| = \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\int_{\Omega} (\mathbf{div} \mathbf{y} + \mathbf{f}) \mathbf{w} \, dx}{\|\nabla \mathbf{w}\|} \leq \mathbf{C}_{\Omega} \|\mathbf{div} \mathbf{y} + \mathbf{f}\|,$$

where \mathbf{C}_{Ω} is the constant in the Friederichs–Steklov inequality for the domain Ω . We observe that

a "noncomputable" negative norm has been estimated by a "computable" one without an attraction of Galerkin orthogonality and local (mesh-dependent) constants.

Thus, for any $\mathbf{y} \in \mathbf{H}(\Omega, \mathbf{div})$ we obtain

$$\begin{aligned} \|\nabla(\mathbf{v} - \mathbf{u})\| &\leq \|\nabla\mathbf{v} - \mathbf{y}\| + \inf_{\mathbf{q} \in \mathbf{Q}_f} \|\mathbf{y} - \mathbf{q}\| \leq \\ &\|\nabla\mathbf{v} - \mathbf{y}\| + \mathbf{C}_\Omega \|\mathbf{div} \mathbf{v} + \mathbf{f}\|. \end{aligned}$$

Above presented *modus operandi* can be viewed as a simplest version of the variational approach to the derivation of Functional Error Majorants.

Deriving a posteriori estimates from integral identities.

For many problems the variational techniques cannot be applied (e.g., because they may have no variational formulation).

In

S. Repin. Two-sided estimates for deviation from an exact solution to uniformly elliptic equation. *Trudi St.-Petersburg Math. Society*, 9(2001), translated in *American Mathematical Translations Series 2*, 9(2003)

it was suggested another method, which is *based on certain transformations of integral identities*.

Non-variational method in the simplest case

Let us expose its simplest version adapted to our model problem.
We have

$$\int_{\Omega} \nabla(\mathbf{u} - \mathbf{v}) \nabla \mathbf{w} \, d\mathbf{x} = \int_{\Omega} (\mathbf{f}\mathbf{w} - \nabla \mathbf{v} \cdot \nabla \mathbf{w}) \, d\mathbf{x}$$

In order to get an upper bound of $\|\nabla(\mathbf{u} - \mathbf{v})\|$ we use the relation

$$\int_{\Omega} (\operatorname{div} \mathbf{y} \mathbf{w} + \nabla \mathbf{w} \cdot \mathbf{y}) \, d\mathbf{x} = 0 \quad \forall \mathbf{w} \in \mathbf{V}_0$$

valid for any $\mathbf{y} \in \mathbf{H}(\Omega, \operatorname{div})$.

We have

$$\begin{aligned}
 & \int_{\Omega} (\nabla \mathbf{v} \cdot \nabla \mathbf{w} - \mathbf{f} \mathbf{w}) \, \mathbf{d}\mathbf{x} = \\
 & \int_{\Omega} (\nabla \mathbf{v} \cdot \nabla \mathbf{w} - \mathbf{f} \mathbf{w} - (\mathbf{div} \mathbf{w} + \nabla \mathbf{w} \cdot \mathbf{y})) \, \mathbf{d}\mathbf{x} = \\
 & \int_{\Omega} ((\nabla \mathbf{v} - \mathbf{y}) \cdot \nabla \mathbf{w} - (\mathbf{f} + \mathbf{div} \mathbf{w}) \mathbf{w}) \, \mathbf{d}\mathbf{x} \leq \\
 & \|\nabla \mathbf{v} - \mathbf{y}\| \|\nabla \mathbf{w}\| + \|\mathbf{f} + \mathbf{div} \mathbf{w}\| \|\mathbf{w}\| \leq \\
 & \leq (\|\nabla \mathbf{v} - \mathbf{y}\| + \mathbf{C}_{\Omega} \|\mathbf{f} + \mathbf{div} \mathbf{w}\|) \|\nabla \mathbf{w}\|.
 \end{aligned}$$

Set $\mathbf{w} = \mathbf{u} - \mathbf{v}$.

$$\int_{\Omega} |\nabla(\mathbf{u} - \mathbf{v})|^2 \mathbf{d}\mathbf{x} \leq (\|\nabla\mathbf{v} - \mathbf{y}\| + \mathbf{C}_{\Omega}\|\mathbf{f} + \mathbf{div}\mathbf{y}\|)\|\nabla(\mathbf{u} - \mathbf{v})\|.$$

Thus, we find that

$$\|\nabla(\mathbf{u} - \mathbf{v})\| \leq \|\nabla\mathbf{v} - \mathbf{y}\| + \mathbf{C}_{\Omega}\|\mathbf{f} + \mathbf{div}\mathbf{y}\|.$$

Functional error estimate. Meaning and properties

For the problem

$$\Delta \mathbf{u} + \mathbf{f} = \mathbf{0}, \quad \mathbf{u} = \mathbf{0} \text{ on } \partial\Omega$$

we have obtained the estimate

$$\|\nabla(\mathbf{u} - \mathbf{v})\| \leq \|\nabla \mathbf{v} - \mathbf{y}\| + C_{\Omega} \|\operatorname{div} \mathbf{y} + \mathbf{f}\| \quad (63)$$

The estimate is valid for any $\mathbf{v} \in \mathbf{V}_0$ and $\mathbf{y} \in \mathbf{H}(\Omega, \operatorname{div})$

Two terms in the right-hand side have a clear sense: they **present measures of the errors in two basic relations**

$$\mathbf{p} = \nabla \mathbf{u}, \quad \operatorname{div} \mathbf{p} + \mathbf{f} = \mathbf{0} \quad \text{in } \Omega$$

that jointly form the equation.

The estimate is sharp

If set $\mathbf{v} = \mathbf{0}$ and $\mathbf{y} = \mathbf{0}$, we obtain the energy estimate for the generalized solution

$$\|\nabla \mathbf{u}\| \leq \mathbf{C}_\Omega \|\mathbf{f}\|.$$

Therefore, no constant less than \mathbf{C}_Ω can be stated in the second term.

If set $\mathbf{y} = \nabla \mathbf{u}$, then the inequality holds as the equality.

Thus, we see that the estimate (63) is **sharp** in the sense that the multipliers of both terms **cannot be taken smaller** and in the set of admissible \mathbf{y} there **exists a function that makes the inequality hold as equality**.

The estimate as a quadratic functional

By means of the algebraic Young's inequality

$$2\mathbf{ab} \leq \beta\mathbf{a}^2 + \frac{1}{\beta}\mathbf{b}^2, \quad \beta > 0$$

we rewrite this estimate in the form

$$\begin{aligned} \|\nabla(\mathbf{u} - \mathbf{v})\|^2 &\leq \\ &\leq (\mathbf{1} + \beta)\|\nabla\mathbf{v} - \mathbf{y}\|^2 + \frac{\mathbf{1} + \beta}{\beta}\mathbf{C}_\Omega^2\|\mathbf{div}\mathbf{y} + \mathbf{f}\|^2 \end{aligned} \quad (64)$$

For any β the right-hand side is a quadratic functional. This property makes it possible to apply well known methods for the minimization with respect to \mathbf{y} .

Deviation Majorant

Denote the right-hand side of (64) by \mathcal{M}_{\oplus} , i.e.,

$$\mathcal{M}_{\oplus}(\mathbf{v}, \mathbf{y}, \beta, \mathbf{C}_{\Omega}, \mathbf{f}) := (\mathbf{1} + \beta) \|\nabla \mathbf{v} - \mathbf{y}\|^2 + \frac{\mathbf{1} + \beta}{\beta} \mathbf{C}_{\Omega}^2 \|\mathbf{div} \mathbf{v} + \mathbf{f}\|^2.$$

This functional provides an upper bound for the norm of the deviation of \mathbf{v} from \mathbf{u} . Therefore, it is natural to call it the **Deviation Majorant**.

BVP $\Delta u + f = 0$ has another variational formulation

$$\begin{aligned} \inf_{\mathbf{v} \in \mathbf{V}_0,} \quad & \mathcal{M}_{\oplus}(\mathbf{v}, \mathbf{y}, \beta, \mathbf{C}_{\Omega}, \mathbf{f}) \\ & \beta > 0, \\ & \mathbf{y} \in \mathbf{H}(\Omega, \text{div}), \end{aligned}$$

- Minimum of this functional is *zero*;
- it is attained if and only if $\mathbf{v} = \mathbf{u}$ and $\mathbf{y} = \mathbf{A}\nabla\mathbf{u}$!;
- \mathcal{M}_{\oplus} contains only one global constant \mathbf{C}_{Ω} , which is problem independent;

In principle, one can select certain sequences of subspaces $\{\mathbf{V}_{hk}\} \in \mathbf{V}_0$ and $\{\mathbf{Y}_{hk}\} \in \mathbf{H}(\Omega, \text{div})$ and minimize the Error Majorant with respect to these subspaces

$$\begin{aligned} & \inf_{\mathbf{v} \in \mathbf{V}_{hk},} \mathcal{M}_{\oplus}(\mathbf{v}, \mathbf{y}, \beta, \mathbf{C}_{\Omega}, \mathbf{f}) \\ & \beta > 0, \\ & \mathbf{y} \in \mathbf{Y}_{hk}, \end{aligned}$$

If the subspaces are limit dense, then we would obtain a sequence of approximate solutions $(\mathbf{v}_k, \mathbf{y}_k)$ and the sequence of numbers

$$\gamma_k := \inf_{\beta > 0} \mathcal{M}_{\oplus}(\mathbf{v}_k, \mathbf{y}_k, \beta, \mathbf{C}_{\Omega}, \mathbf{f}) \rightarrow 0$$

How to use the Majorants in practice?

Consider **CONFORMING FEM APPROXIMATIONS**.

We have 3 basic ways to use the deviation estimate:

- (a) **Direct** (via flux averaging on the mesh \mathcal{T}_h);
- (b) **One step delay** (via flux averaging on the mesh \mathbf{h}_{ref});
- (c) **Minimization** (minimization via \mathbf{y}).

(a) Use recovered gradients

Let $\mathbf{u}_h \in \mathbf{V}_h$, then

$$\mathbf{p}_h := \nabla \mathbf{u}_h \in \mathbf{L}_2(\Omega, \mathbb{R}^d), \quad \mathbf{p}_h \notin \mathbf{H}(\Omega, \text{div}).$$

Use an averaging operator $\mathbf{G}_h : \mathbf{L}_2(\Omega, \mathbb{R}^d) \rightarrow \mathbf{H}(\Omega, \text{div})$ and have a **directly computable estimate**

$$\|\nabla(\mathbf{u} - \mathbf{u}_h)\| \leq \|\nabla \mathbf{u}_h - \mathbf{G}_h \mathbf{p}_h\| + \mathbf{C}_\Omega \|\text{div} \mathbf{G}_h \mathbf{p}_h + \mathbf{f}\|$$

(b) Use recovered gradients from $\mathcal{T}_{h_{ref}}$

Let $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k, \dots$ be a sequence of approximations on meshes \mathcal{T}_{h_k} .
Compute $\mathbf{p}_k := \nabla \mathbf{u}_k$, average it by \mathbf{G}_k and for \mathbf{u}_{k-1} use the estimate

$$\|\mathbf{u} - \mathbf{u}_{k-1}\| \leq \|\nabla \mathbf{u}_{k-1} - \mathbf{G}_k \mathbf{p}_k\| + \mathbf{C}_\Omega \|\operatorname{div} \mathbf{G}_k \mathbf{p}_k + \mathbf{f}\|$$

This estimate gives **a quantitative form of the Runge's rule.**

(c) Minimize \mathcal{M}_{\oplus} with respect to y .

Select a certain subspace \mathbf{Y}_{τ} in $\mathbf{H}(\Omega, \mathbf{div})$. **Generally, \mathbf{Y}_{τ} may be constructed on another mesh \mathcal{T}_{τ} and with help of different trial functions.**

Then

$$\|\nabla(\mathbf{u} - \mathbf{u}_h)\| \leq \inf_{\mathbf{y}_h \in \mathbf{Y}_h} \{ \|\nabla \mathbf{u}_h - \mathbf{y}_h\| + \mathbf{C}_{\Omega} \|\mathbf{div} \mathbf{y}_h + \mathbf{f}\| \}$$

The wider $\mathbf{Y}_h \subset \mathbf{H}(\Omega, \mathbf{div})$ the sharper is the upper bound.

Quadratic type functional

From the technical point of view it is better to square both parts of the estimate and apply minimization to a quadratic functional, namely

$$\begin{aligned} \|\nabla(\mathbf{u} - \mathbf{u}_h)\|^2 &\leq \inf_{\mathbf{y}_h \in \mathbf{Y}_h} \left\{ (\mathbf{1} + \beta) \|\nabla \mathbf{u}_h - \mathbf{y}_h\| + \right. \\ &\quad \left. + \mathbf{C}_\Omega \left(\mathbf{1} + \frac{\mathbf{1}}{\beta} \right) \|\mathbf{div} \mathbf{y}_h + \mathbf{f}\|^2 \right\} \end{aligned}$$

Here, the positive parameter β can be also used to minimize the right-hand side.

Before going to more complicated problems where Deviation Majorants are derived by a more sophisticated theory, we observe several simple examples that nevertheless reflect key points of the above method.

Simple 1-D problem

$$\begin{aligned}(\alpha(x) u')' &= f(x), \\ u(a) &= 0, \quad u(b) = u_b.\end{aligned}$$

It is equivalent to the variational problem

$$J(v) = \int_a^b \left(\frac{1}{2} \alpha(x) |v'|^2 + f(x)v \right) dx.$$

Assume that the coefficient α belongs to $\in \mathbf{L}^\infty$ and bounded from below by a positive constant. Now

$$\mathbf{V}_0 + \mathbf{u}_0 = \{v \in \mathbf{H}^1(a, b) \mid v(a) = 0, v(b) = u_b\}.$$

Deviation Majorant

$$\mathcal{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y}) = (\mathbf{1} + \beta) \left(\int_a^b |\alpha \mathbf{v}' - \mathbf{y}|^2 \, d\mathbf{x} + \frac{\mathbf{C}_{(a,b)}^2}{\beta} \int_a^b |\mathbf{y}' - \mathbf{f}|^2 \, d\mathbf{x} \right). \quad (65)$$

In this simple model, \mathbf{u} can be presented in the form

$$\mathbf{u}(\mathbf{x}) = \int_a^{\mathbf{x}} \frac{\mathbf{1}}{\alpha(\mathbf{t})} \int_a^{\mathbf{t}} \mathbf{f}(\mathbf{z}) \, d\mathbf{z} \, d\mathbf{t} + \frac{\mathbf{x}}{\mathbf{b}} \left(\mathbf{u}_b - \int_a^{\mathbf{b}} \frac{\mathbf{1}}{\alpha(\mathbf{t})} \int_a^{\mathbf{t}} \mathbf{f}(\mathbf{z}) \, d\mathbf{z} \, d\mathbf{t} \right).$$

what gives an opportunity to verify how error estimation methods work.

Approximations

Let \mathbf{V}_h be made of piecewise- \mathbf{P}^1 continuous functions on uniform splittings of the interval and consider approximations of the following types:

- Galerkin approximations;
- Approximations very close to Galerkin (sharp);
- Approximations which are "good" but not Galerkin;
- Coarse (rough) approximations.

Our aim is to show that the Deviation Majorant can be effectively used as an error estimation instrument in all the above cases.

Computation of the Majorant

To find a sharp upper bound, we minimize \mathcal{M}_{\oplus} with respect to \mathbf{y} and β starting from the function $\mathbf{y}_0 = \mathbf{G}(\mathbf{v}')$, where \mathbf{G} is a simple averaging operator, e.g, defined by the relations

$$\mathbf{G}(\mathbf{v}')(x_i) = \frac{1}{2}(\mathbf{v}'(x_i - \mathbf{0}) + \mathbf{v}'(x_i + \mathbf{0})),$$

By the quantity

$$\inf_{\beta > 0} \mathcal{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y}_0),$$

we obtain a coarse upper bound of the error. It is further improved by minimizing \mathcal{M}_{\oplus} with respect to \mathbf{y} .

Example

Let $\alpha(x) = 1$, $f(x) = c$, $a = 0$, $b = 1$, $u_b = 1$, e.g., we consider the problem

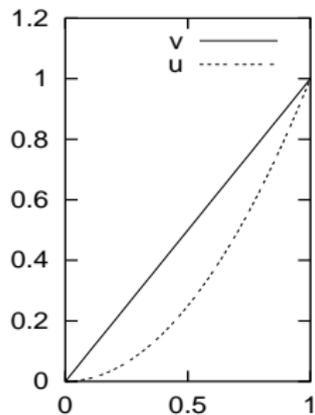
$$u'' = 2, \quad u(0) = 0, \quad u(1) = 1.$$

In this case, $C_{(a,b)} = 1/\pi$ and

$$u = \frac{c}{2}x^2 + \left(1 - \frac{c}{2}\right)x, \quad u' = cx + 1 - \frac{c}{2}.$$

Take a **rough** approximation $v = x$. Then

$$\|(u - v)'\|^2 = \int_0^1 c^2(x - 0.5)^2 dx = c^2/12 \approx 0.083c^2.$$



Exact solution and an approximation.

Various \mathbf{y} give different upper bounds

(a) Take $\mathbf{y} = \mathbf{v}' = \mathbf{1}$, then the first term in

$$\mathcal{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y}) = (\mathbf{1} + \beta) \left(\int_0^1 |\mathbf{v}' - \mathbf{y}|^2 \, dx + \frac{\mathbf{1}}{\pi^2 \beta} \int_0^1 |\mathbf{y}' - \mathbf{f}|^2 \, dx \right) dx.$$

vanishes and we have $\mathcal{M}_{\oplus} \rightarrow \mathbf{c}^2/\pi^2 \approx \mathbf{0.101c}^2$; as $\beta \rightarrow +\infty$. We see that this upper bound overestimates true error. Note that in this case, all sensible averagings of $\mathbf{v}' = \mathbf{1}$ give exactly the same function: $\mathbf{G}(\mathbf{1}) = \mathbf{1}!$ Therefore,

$$\mathbf{G}(\mathbf{v}') - \mathbf{v}' \equiv \mathbf{0}$$

and formally ZZ indicator "does not see the error".

For the choice $\mathbf{y} = \mathbf{v}'$ the Majorant give a certain upper bound of the error (which is not so bad), but the integrand cannot indicate the distribution of local errors. Indeed, we have

$$\mathcal{M}_{\oplus} = \frac{1}{\pi^2} \int_0^1 \mathbf{c}^2 \mathbf{d}\mathbf{x}.$$

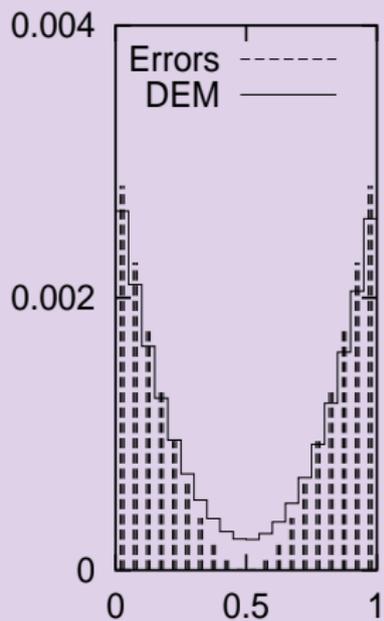
However, the integrand of the Majorant is a constant function, but the error is distributed in accordance with a parabolic law:

$$(\mathbf{u} - \mathbf{v})' = \mathbf{c}(\mathbf{x} - 0.5)^2.$$

(b). Take $\mathbf{y} = \mathbf{c}\mathbf{x} + \mathbf{1} - \mathbf{c}/2$. Then, $\mathbf{y}' = \mathbf{c}$ and the second term of the majorant vanishes. We have (for $\beta = 0$)

$$\mathcal{M}_{\oplus} = \int_0^1 \mathbf{c}^2 (\mathbf{x} - 1/2)^2 d\mathbf{x} = \mathbf{c}^2/12.$$

We observe that both the global error and the error distribution are exactly reproduced. In real life computations such an "ideal" function \mathbf{y} may be unattainable. However, the minimization makes the Majorant close to the sharp value. In this elementary example, we have minimized the Majorant on using piecewise affine approximations of \mathbf{y} on 20 subintervals. The elementwise error distribution obtained as the result of this procedure is exposed on the next picture.



True errors and those computed by the Majorant.

To give further illustrations, we consider the functions

$$\mathbf{u}_\delta = \mathbf{u} + \delta\varphi,$$

where δ is a number and φ is a certain function (perturbation).

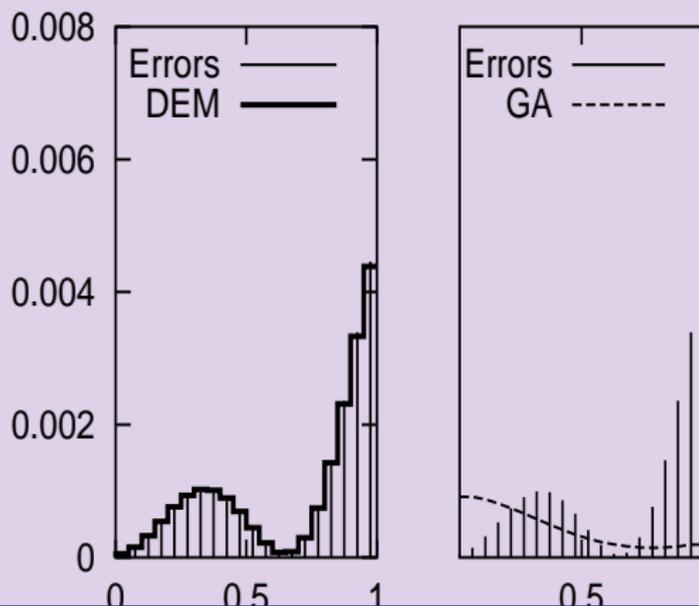
Approximate solutions (whose errors are measured) are piecewise affine continuous interpolants of \mathbf{u}_δ defined on a uniform mesh with 20 subintervals. We take $\varphi = \mathbf{x} \sin(\pi\mathbf{x})$ and $\delta = 0.1, 0.01, 0.001$, and 0.

Table:

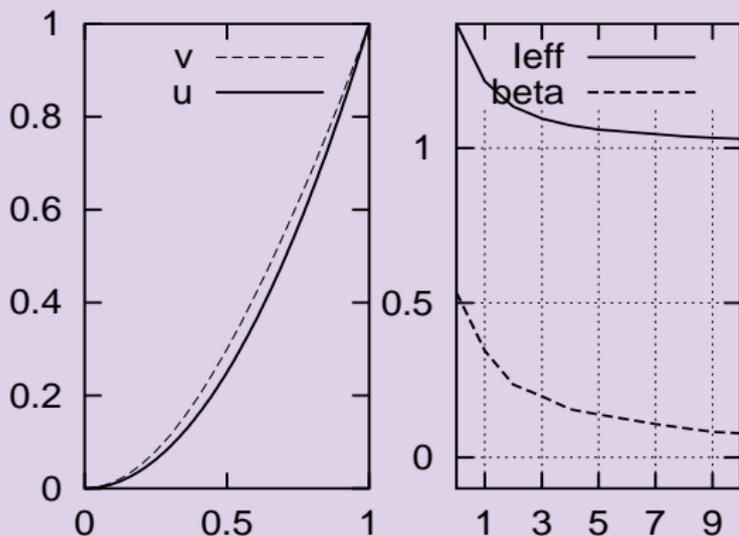
δ	e	$2\mathcal{M}_\oplus$	$2\mathcal{M}_\ominus$	i_{eff}	i_{esh}
0.1	0.019692	0.019743	0.019683	1.003	1.018
0.01	0.001022	0.001025	0.001013	1.003	1.011
0.001	0.000835	0.000839	0.000827	1.005	1.002
0	0.000833	0.000836	0.000825	1.004	1.002

In this experiment the Majorant was computed for $\frac{1}{2}\|\mathbf{e}\|^2$.

Error estimation for $\delta = 0.1$

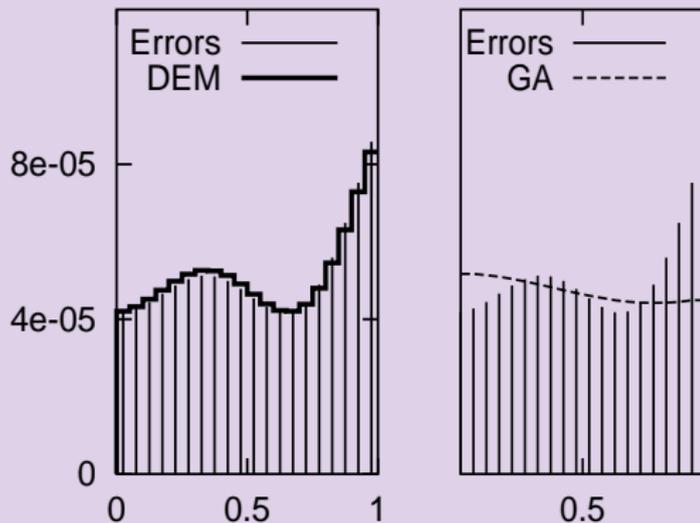


Functions \mathbf{u} , \mathbf{v} and \mathbf{i}_{eff} for $\delta = 0.1$

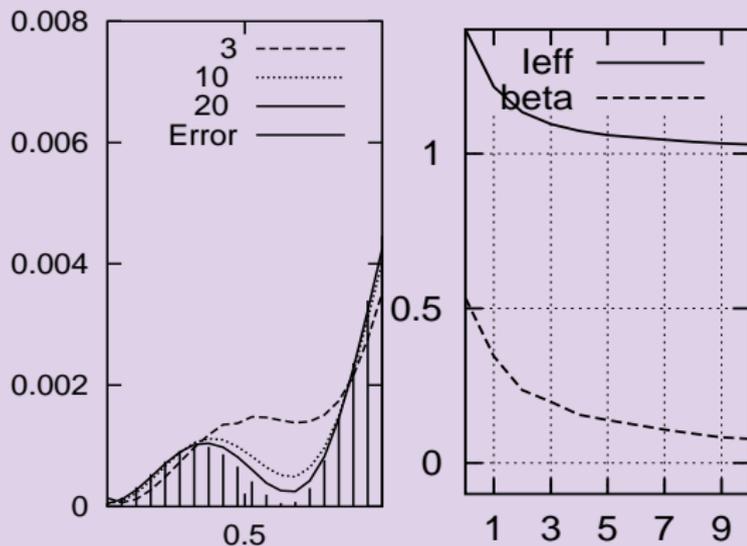


Error estimation for $\delta = 0.01$

A more precise approximation.

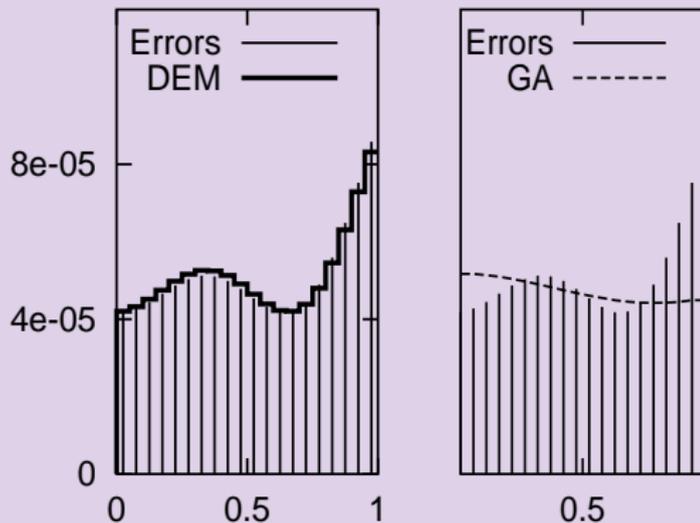


Functions $\mathbf{e}(\mathbf{y})$, β and \mathbf{i}_{eff} for $\delta = 0.1$

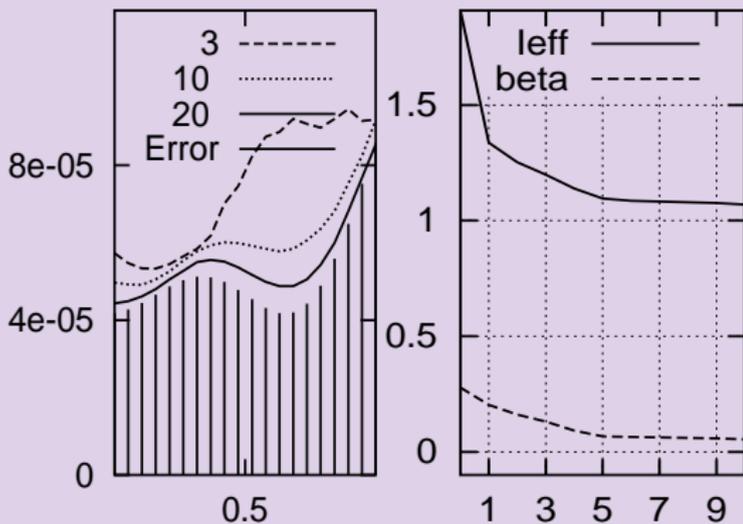


Error estimation for $\delta = 0.01$

A more precise approximation.

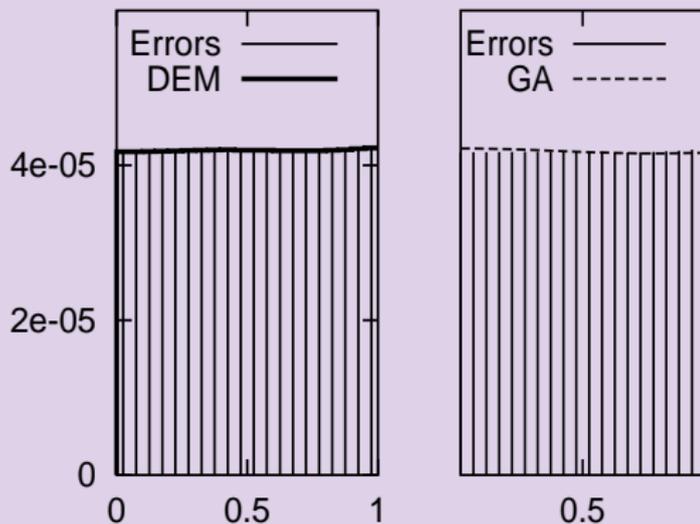


Functions $\mathbf{e}(\mathbf{y})$, β and \mathbf{i}_{eff} for $\delta = 0.01$

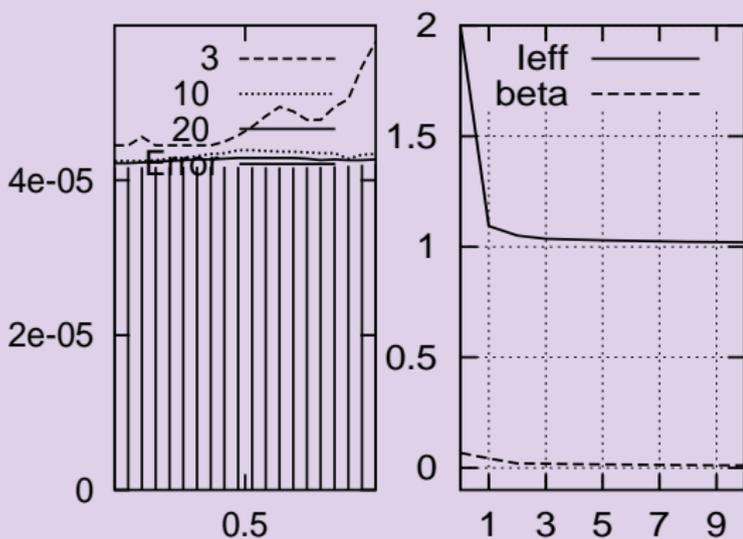


Error estimation for $\delta = 0.001$

Sharp approximation.

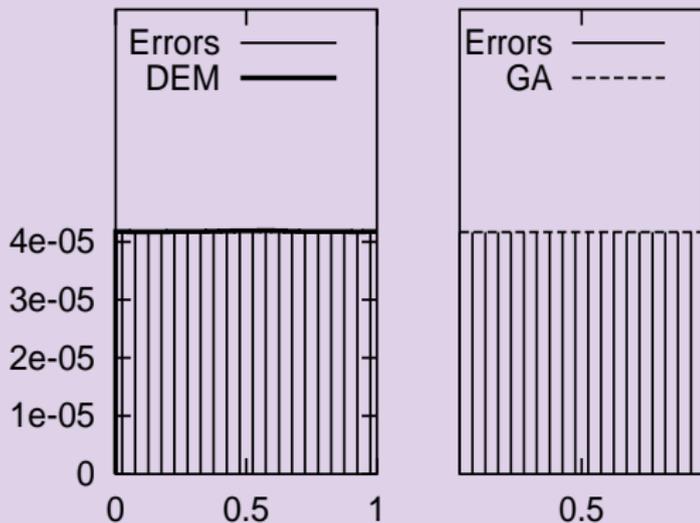


Functions $\mathbf{e}(\mathbf{y})$, β and \mathbf{i}_{eff} for $\delta = 0.001$

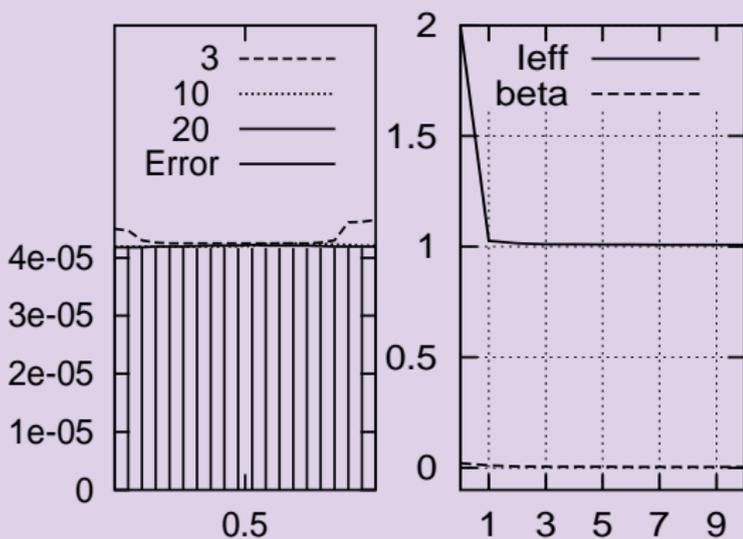


Error estimation for $\delta = 0$

Interpolant of the exact solution.



Functions $\mathbf{e}(\mathbf{y})$, β and \mathbf{i}_{eff} for $\delta = 0$



Task 6

Apply the above theory to the problem

$$(\alpha \mathbf{u}')' = \mathbf{f},$$

$$\mathbf{u}(0) = \mathbf{0}, \quad \mathbf{u}(1) = \mathbf{b}$$

with your own α , \mathbf{f} , and \mathbf{b} . Compute approximate solutions and verify their accuracy along the same lines as in the example above.

Other examples

For problems with lower terms it is easy to obtain estimates without C_Ω .

$$\begin{aligned}\Delta \mathbf{u} - \varrho \mathbf{u} + \mathbf{f} &= \mathbf{0}, & \varrho > 0, \\ \mathbf{u} &= \mathbf{u}_0 \quad \text{on } \partial\Omega.\end{aligned}$$

Such estimates can be derived by both *variational* and *non-variational* method. Let $\mathbf{w} \in \mathbf{V}_0 := \mathring{\mathbf{H}}^1(\Omega)$. We have

$$\begin{aligned}& \int_{\Omega} \nabla(\mathbf{u} - \mathbf{v}) \cdot \nabla \mathbf{w} \, \mathbf{d}\mathbf{x} + \varrho \int_{\Omega} (\mathbf{u} - \mathbf{v}) \mathbf{w} \, \mathbf{d}\mathbf{x} = \\ &= \int_{\Omega} (\mathbf{f}\mathbf{w} - \nabla \mathbf{v} \cdot \nabla \mathbf{w}) \, \mathbf{d}\mathbf{x} - \varrho \int_{\Omega} \mathbf{v}\mathbf{w} \, \mathbf{d}\mathbf{x}.\end{aligned}$$

Use the integral identity for $\mathbf{y} \in \mathbf{H}(\Omega, \text{div})$:

$$\int_{\Omega} (\nabla \mathbf{w} \cdot \mathbf{y} + \mathbf{w} \text{div} \mathbf{y}) \, dx = 0 \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

$$\begin{aligned} \int_{\Omega} \nabla(\mathbf{u} - \mathbf{v}) \cdot \nabla \mathbf{w} \, dx + \varrho \int_{\Omega} (\mathbf{u} - \mathbf{v}) \mathbf{w} \, dx &= \\ \int_{\Omega} (\mathbf{f} + \text{div} \mathbf{y} - \varrho \mathbf{v}) \mathbf{w} \, dx + \int_{\Omega} (\mathbf{y} - \nabla \mathbf{v}) \cdot \nabla \mathbf{w} \, dx &\leq \\ \leq \|\mathbf{f} + \text{div} \mathbf{y} - \varrho \mathbf{v}\| \|\mathbf{w}\| + \|\nabla \mathbf{v} - \mathbf{y}\| \|\nabla \mathbf{w}\|. \end{aligned}$$

Set $\mathbf{w} = \mathbf{u} - \mathbf{v}$ and note that

$$\begin{aligned} & \| \mathbf{f} + \mathbf{div} \boldsymbol{\rho} \mathbf{v} \| \| \mathbf{u} - \mathbf{v} \| + \| \nabla \mathbf{v} - \mathbf{y} \| \| \nabla (\mathbf{u} - \mathbf{v}) \| = \\ &= \frac{1}{\varrho} \| \mathbf{f} + \mathbf{div} \boldsymbol{\rho} \mathbf{v} \| \varrho \| \mathbf{u} - \mathbf{v} \| + \| \nabla \mathbf{v} - \mathbf{y} \| \| \nabla (\mathbf{u} - \mathbf{v}) \| \leq \\ &\leq \left(\frac{1}{\varrho^2} \| \mathbf{f} + \mathbf{div} \boldsymbol{\rho} \mathbf{v} \|^2 + \| \nabla \mathbf{v} - \mathbf{y} \|^2 \right)^{1/2} \| \mathbf{u} - \mathbf{v} \| \end{aligned}$$

where

$$\| \mathbf{u} - \mathbf{v} \|^2 = \int_{\Omega} (|\nabla (\mathbf{u} - \mathbf{v})|^2 + \varrho |\mathbf{u} - \mathbf{v}|^2) \mathbf{d}x.$$

Then, we obtain the estimate

$$\| \mathbf{u} - \mathbf{v} \|^2 \leq \frac{1}{\varrho^2} \| \mathbf{f} + \mathbf{div} \mathbf{v} - \varrho \mathbf{v} \|^2 + \| \nabla \mathbf{v} - \mathbf{y} \|^2$$

By the variational method this estimate was derived in 97'. Also, it readily follows from the general a posteriori framework (see, e.g., S.R. Math. Comp. 2000).

This estimate has no \mathbf{C}_Ω . However, in practice, it may give *big overestimation* if ϱ is small due to large penalty at the first term.

How Functional A Posteriori Estimate looks like for the problem

$$\mathbf{div} \mathbf{A} \nabla \mathbf{u} + \mathbf{f} = \mathbf{0}, \quad \mathbf{u} = \mathbf{u}_0 \quad \text{on } \partial\Omega$$

and for problems with boundary conditions of other types?

For the generalized diffusion equation it is as follows:

$$\| \nabla(\mathbf{u} - \mathbf{v}) \|^2 \leq \sqrt{\mathbf{D}(\nabla \mathbf{v}, \mathbf{y})} + \mathbf{C}(\Omega, \mathbf{A}) \|\mathbf{div} \mathbf{y} + \mathbf{f}\|,$$

where

$$\begin{aligned} \|\eta\|^2 &:= \int_{\Omega} \mathbf{A} \nabla \eta \cdot \nabla \eta \, \mathbf{d}\mathbf{x}, \\ \mathbf{D}(\nabla \mathbf{v}, \mathbf{y}) &:= \int_{\Omega} (\mathbf{A} \nabla \mathbf{v} \cdot \nabla \mathbf{v} + \mathbf{A}^{-1} \mathbf{y} \cdot \mathbf{y} - 2 \nabla \mathbf{v} \cdot \mathbf{y}) \, \mathbf{d}\mathbf{x} = \\ &= \int_{\Omega} (\mathbf{A} \nabla \mathbf{v} - \mathbf{y}) \cdot (\nabla \mathbf{v} - \mathbf{A}^{-1} \mathbf{y}) \, \mathbf{d}\mathbf{x}, \\ \|\mathbf{w}\| &\leq \mathbf{C}(\Omega, \mathbf{A}) \|\nabla \mathbf{w}\|. \end{aligned}$$

How such estimates are derived we will discuss in the next Lecture.

Lecture 4

We expose the general approach to deriving **two-sided functional estimates of the deviations** from exact solutions of linear elliptic type problems having the operator form

$$\Lambda^* \mathcal{A} \Lambda u + \ell = 0$$

where Λ and \mathcal{A} are linear bounded operators and \mathcal{A} is symmetric and positive definite.

Lecture plan

- Two-sided a posteriori estimates for linear elliptic type problems;
- Properties: computability, consistency, reliability;
- Relationships with other error estimation methods;
- Diffusion equation with Dirichlet boundary conditions;
- Diffusion equation with Neumann boundary conditions;
- Diffusion equation with mixed boundary conditions;
- Linear elasticity with mixed boundary conditions;

Problem in the abstract form

Many problems can be presented in the following form: **find**
 $\mathbf{u} \in \mathbf{V}_0 + \mathbf{u}_0$ such that

$$(\mathcal{A}\Lambda\mathbf{u}, \Lambda\mathbf{w}) + \langle \boldsymbol{\ell}, \mathbf{w} \rangle = 0 \quad \forall \mathbf{w} \in \mathbf{V}_0. \quad (66)$$

Here \mathbf{V}_0 is a subspace of a reflexive Banach space \mathbf{V} ,

e.g., $\mathbf{V} = \mathbf{H}^1$, $\mathbf{V}_0 = \overset{\circ}{\mathbf{H}}^1$.

$\Lambda : \mathbf{V} \rightarrow \mathbf{U}$ is a bounded linear operator, e.g. $\Lambda = \nabla$.

\mathbf{U} is a Hilbert space with scalar product (\cdot, \cdot) and norm $\|\cdot\|$,
 e.g., $\mathbf{U} = \mathbf{L}^2$.

$\boldsymbol{\ell} \in \mathbf{V}_0^*$, is a linear functional in the dual space, e.g., in \mathbf{H}^{-1} . In
 general, we may assume that

$$\langle \boldsymbol{\ell}, \mathbf{w} \rangle = (\mathbf{f}, \mathbf{w}) + (\mathbf{g}, \Lambda\mathbf{w}).$$

$\mathcal{A} \in \mathcal{L}(\mathbf{U}, \mathbf{U})$ is a self-adjoint operator.

Assumptions

We assume that

$$\mathbf{V} \text{ is compactly embedded in } \mathbf{U} \quad (67)$$

and the operators \mathbf{A} and \mathcal{A} satisfy the relations

$$\mathbf{c}_1 \|\mathbf{y}\|^2 \leq (\mathcal{A}\mathbf{y}, \mathbf{y}) \leq \mathbf{c}_2 \|\mathbf{y}\|^2, \quad \forall \mathbf{y} \in \mathbf{U}, \quad (68)$$

$$\|\mathbf{A}\mathbf{w}\| \geq \mathbf{c}_3 \|\mathbf{w}\|_{\mathbf{V}}, \quad \forall \mathbf{w} \in \mathbf{V}_0, \quad (69)$$

For our analysis, it is convenient to introduce two more norms:

$$\| \mathbf{y} \| := (\mathcal{A}\mathbf{y}, \mathbf{y})^{1/2}, \quad \| \mathbf{y} \|_* := (\mathcal{A}^{-1}\mathbf{y}, \mathbf{y})^{1/2},$$

where $\mathcal{A}^{-1} : \mathbf{U} \rightarrow \mathbf{U}$

is the operator inverse to \mathcal{A} . The respective spaces \mathbf{Y} and \mathbf{Y}^* contain elements of \mathbf{U} equipped with the norms $\| \cdot \|$ and $\| \cdot \|_*$, respectively. Let $\mathbf{\Lambda}^*$ be the operator conjugate to $\mathbf{\Lambda}$, i.e.,

$$(\mathbf{y}, \mathbf{\Lambda}\mathbf{w}) = \langle \mathbf{\Lambda}^*\mathbf{y}, \mathbf{w} \rangle, \quad \forall \mathbf{w} \in \mathbf{V}_0. \quad (70)$$

Getting a posteriori estimates by transformations of integral identities

For a detailed exposition see

S.Repin. Two-sided estimates of deviation from exact solutions of uniformly elliptic equations, Proc. St. Petersburg Math. Society, IX(2001), pp. 143–171, translation in Amer. Math. Soc. Transl. Ser. 2, 209, Amer. Math. Soc., Providence, RI, 2003.

Let $\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0$ be an approximation.

$$(\mathcal{A}\Lambda(\mathbf{u} - \mathbf{v}), \Lambda\mathbf{w}) + \langle \ell, \mathbf{w} \rangle + (\mathcal{A}\Lambda\mathbf{v}, \Lambda\mathbf{w}) = 0 \quad \forall \mathbf{w} \in \mathbf{V}_0,$$

Set $\mathbf{w} = \mathbf{u} - \mathbf{v}$.

$$\| \Lambda(\mathbf{u} - \mathbf{v}) \|^2 = | \langle \ell, \mathbf{u} - \mathbf{v} \rangle + (\mathcal{A}\Lambda\mathbf{v}, \Lambda(\mathbf{u} - \mathbf{v})) | .$$

By

$$(\mathbf{y}, \Lambda\mathbf{w}) = \langle \Lambda^*\mathbf{y}, \mathbf{w} \rangle, \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

we have

$$\| \Lambda(\mathbf{u} - \mathbf{v}) \|^2 = | \langle \ell, \mathbf{u} - \mathbf{v} \rangle + (\mathcal{A}\Lambda\mathbf{v} - \mathbf{y}, \Lambda(\mathbf{u} - \mathbf{v})) + \langle \Lambda^*\mathbf{y}, \mathbf{u} - \mathbf{v} \rangle | .$$

Therefore, we find that

$$\begin{aligned} \|\Lambda(\mathbf{u} - \mathbf{v})\|^2 &\leq | \langle \ell + \Lambda^* \mathbf{y}, \mathbf{u} - \mathbf{v} \rangle | + | (\mathcal{A}\Lambda\mathbf{v} - \mathbf{y}, \Lambda(\mathbf{u} - \mathbf{v})) | \leq \\ &\leq \|\ell + \Lambda^* \mathbf{y}\| \|\Lambda(\mathbf{u} - \mathbf{v})\| + \|\mathcal{A}\Lambda\mathbf{v} - \mathbf{y}\|_* \|\Lambda(\mathbf{u} - \mathbf{v})\|, \end{aligned}$$

where

$$\|\mu\| := \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\langle \mu, \mathbf{w} \rangle}{\|\Lambda\mathbf{w}\|}$$

denotes the norm of the functional $\mu : \mathbf{V}_0 \rightarrow \mathbb{R}$.

To prove that the value of $\|\ell + \Lambda^* \mathbf{y}\|$ is finite we note that

$$\begin{aligned} \langle \ell + \Lambda^* \mathbf{y}, \mathbf{w} \rangle &\leq \|\ell\|_{\mathbf{V}_0^*} \|\mathbf{w}\|_{\mathbf{V}} + \|\mathbf{y}\| \|\Lambda\mathbf{w}\| \leq (\mathbf{c}_3^{-1} \|\ell\|_{\mathbf{V}_0^*} + \|\mathbf{y}\|) \|\Lambda\mathbf{w}\| \leq \\ &\leq \mathbf{c}_1^{-1/2} (\mathbf{c}_3^{-1} \|\ell\|_{\mathbf{V}_0^*} + \|\mathbf{y}\|) \|\Lambda\mathbf{w}\|. \end{aligned}$$

General estimate

As a result we obtain the general form of a functional a posteriori estimate for an elliptic type problem:

$$\| \Lambda(\mathbf{u} - \mathbf{v}) \| \leq \mathbf{I} \ell + \Lambda^* \mathbf{y} \mathbf{I} + \| \mathcal{A}\Lambda\mathbf{v} - \mathbf{y} \|_* . \quad (71)$$

Denote

$$\begin{aligned}
 \mathbf{D}(\mathbf{y}_1, \mathbf{y}_2) &:= \frac{1}{2}(\mathcal{A}\mathbf{y}_1, \mathbf{y}_1) + \frac{1}{2}(\mathcal{A}^{-1}\mathbf{y}_2, \mathbf{y}_2) - (\mathbf{y}_1, \mathbf{y}_2) = \\
 &= \frac{1}{2}(\mathcal{A}(\mathbf{y}_1 - \mathcal{A}^{-1}\mathbf{y}_2), \mathbf{y}_1 - \mathcal{A}^{-1}\mathbf{y}_2) = \frac{1}{2} \|\mathbf{y}_1 - \mathcal{A}^{-1}\mathbf{y}_2\|^2 = \\
 &= \frac{1}{2}(\mathcal{A}^{-1}(\mathbf{y}_2 - \mathcal{A}\mathbf{y}_1), \mathbf{y}_2 - \mathcal{A}\mathbf{y}_1) = \frac{1}{2} \|\mathbf{y}_2 - \mathcal{A}\mathbf{y}_1\|_*^2
 \end{aligned}$$

Then

$$\|\mathcal{A}\boldsymbol{\Lambda}\mathbf{v} - \mathbf{y}\|_*^2 = 2\mathbf{D}(\boldsymbol{\Lambda}\mathbf{v}, \mathbf{y})$$

and we obtain

$$\| \Lambda(\mathbf{u} - \mathbf{v}) \| \leq \| \ell + \Lambda^* \mathbf{y} \| + \sqrt{2\mathbf{D}(\Lambda \mathbf{v}, \mathbf{y})}.$$

Square both sides and use Young's inequality

$$\| \Lambda(\mathbf{u} - \mathbf{v}) \|^2 \leq \left(1 + \frac{1}{\beta}\right) \| \ell + \Lambda^* \mathbf{y} \|^2 + 2(1 + \beta) \mathbf{D}(\Lambda \mathbf{v}, \mathbf{y}). \quad (72)$$

This estimate holds for any $\mathbf{y} \in \mathbf{Y}^*$ and $\beta > 0$. Denote its right-hand side by $2\mathcal{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y})$

$2\mathcal{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y})$ is a sharp upper bound

Proposition

For any $\beta > 0$ there exists $\mathbf{y} \in \mathbf{Y}^*$ such that

$$2\mathcal{M}_{\oplus}(\beta, \mathbf{v}, \mathbf{y}) = \|\Lambda(\mathbf{u} - \mathbf{v})\|^2.$$

Proof. Set $\mathbf{y}_1 = \frac{1}{1+\beta}(\mathbf{p} + \beta\mathcal{A}\Lambda\mathbf{v})$ where $\mathbf{p} = \mathcal{A}\Lambda\mathbf{u}$. Note that $\langle \ell + \Lambda^*\mathbf{y}_1, \mathbf{w} \rangle = \langle -\mathbf{p} + \mathbf{y}_1, \Lambda\mathbf{w} \rangle$.

$$\begin{aligned} \|\ell + \Lambda^*\mathbf{y}_1\| &= \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{(\mathbf{y}_1 - \mathbf{p}, \Lambda\mathbf{w})}{\|\Lambda\mathbf{w}\|} = \\ &= \frac{\beta}{1+\beta} \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{(\mathcal{A}\Lambda(\mathbf{v} - \mathbf{u}), \Lambda\mathbf{w})}{\|\Lambda\mathbf{w}\|} = \frac{\beta}{1+\beta} \|\Lambda(\mathbf{v} - \mathbf{u})\|, \end{aligned}$$

Similarly

$$\mathbf{D}(\Lambda \mathbf{v}, \mathbf{y}_1) = \frac{1}{2(1+\beta)^2} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2 .$$

Therefore,

$$\begin{aligned} \left(1 + \frac{1}{\beta}\right) \|\ell + \Lambda^* \mathbf{y}\|^2 + 2(1 + \beta) \mathbf{D}(\Lambda \mathbf{v}, \mathbf{y}) &= \\ \left(1 + \frac{1}{\beta}\right) \left(\frac{\beta}{1+\beta} \|\Lambda(\mathbf{v} - \mathbf{u})\|\right)^2 + 2(1 + \beta) \left(\frac{1}{2(1+\beta)^2} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2\right) &= \\ = \frac{\beta}{1+\beta} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2 + \frac{1}{1+\beta} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2 = \|\Lambda(\mathbf{v} - \mathbf{u})\|^2 . \end{aligned}$$

We replace $[\ell + \Lambda^* \mathbf{y}]$ by the norm in a Hilbert space \mathbf{U} provided that ℓ belongs to a narrower set. Assume that

$$\begin{aligned} \ell &\in \mathbf{U} \subset \mathbf{V}_0^*, \\ \mathbf{y} &\in \mathbf{Q}^* := \{\mathbf{z}^* \in \mathbf{Y}^* \mid \Lambda^* \mathbf{z}^* \in \mathbf{U}\}. \end{aligned}$$

Note that \mathbf{Q}^* can be endowed with the norm

$$\|\mathbf{y}\|_{\mathbf{Q}^*}^2 := \|\mathbf{y}\|_*^2 + \|\Lambda^* \mathbf{y}\|_{\mathbf{U}}^2.$$

If $\ell \in \mathbf{U}$, then \mathbf{Q}^* contains the exact solution \mathbf{p} of the dual problem! This fact is important for the proof of the sharpness of the Majorant.

Then

$$\langle \ell + \Lambda^* \mathbf{y}, \mathbf{w} \rangle = (\ell + \Lambda^* \mathbf{y}, \mathbf{w}) \quad \mathbf{w} \in \mathbf{V}_0.$$

$$\begin{aligned} \|\ell + \Lambda^* \mathbf{y}\| &= \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\langle \ell + \Lambda^* \mathbf{y}, \mathbf{w} \rangle}{\|\Lambda \mathbf{w}\|} \leq \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\|\ell + \Lambda^* \mathbf{y}\| \|\mathbf{w}\|}{\|\Lambda \mathbf{w}\|} \leq \\ &\leq \|\ell + \Lambda^* \mathbf{y}\| \mathbf{c}_1^{-1} \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\|\mathbf{w}\|}{\|\Lambda \mathbf{w}\|} \leq \mathbf{c}_1^{-1} \mathbf{c}_3^{-1} \|\ell + \Lambda^* \mathbf{y}\|. \end{aligned}$$

Denote $\mathbf{c}^2 = \mathbf{c}_1^{-2} \mathbf{c}_3^{-2}$.

Computable Majorant of the deviation

Now, the Majorant \mathcal{M}_{\oplus} is replaced by \mathbf{M}_{\oplus} , namely we arrive at the estimate

$$\begin{aligned} \frac{1}{2} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2 &\leq \mathbf{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y}) := \\ &:= (\mathbf{1} + \beta) \mathbf{D}(\Lambda \mathbf{v}, \mathbf{y}) + \frac{\mathbf{1} + \beta}{2\beta} \mathbf{c}^2 \|\ell + \Lambda^* \mathbf{y}\|^2. \quad (73) \end{aligned}$$

Variational Method

Problem \mathcal{P} . Find $\mathbf{u} \in \mathbf{V}_0 + \mathbf{u}_0$ such that

$$\mathbf{J}(\mathbf{u}) = \inf_{\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0} \mathbf{J}(\mathbf{u}) := \inf \mathcal{P},$$

where

$$\mathbf{J}(\mathbf{v}) = \frac{1}{2} \|\Lambda \mathbf{v}\|^2 + \langle \ell, \mathbf{v} \rangle.$$

Lagrangian

On the set $(\mathbf{V}_0 + \mathbf{u}_0) \times \mathbf{Y}^*$, we define the Lagrangian

$$\mathbf{L}(\mathbf{v}, \mathbf{y}) = (\mathbf{y}, \Lambda \mathbf{v}) - \frac{1}{2} \|\mathbf{y}\|^2 + \langle \ell, \mathbf{v} \rangle$$

and the functional

$$\mathbf{I}^*(\mathbf{y}) = \inf_{\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0} \mathbf{L}(\mathbf{v}, \mathbf{y}) = \begin{cases} (\mathbf{y}, \Lambda \mathbf{u}_0) - \frac{1}{2} \|\mathbf{y}\|_*^2 + \langle \ell, \mathbf{u}_0 \rangle, & \mathbf{y} \in \mathbf{Q}_\ell^*, \\ -\infty, & \mathbf{y} \notin \mathbf{Q}_\ell^*, \end{cases}$$

where $\mathbf{Q}_\ell^* := \{\mathbf{y} \in \mathbf{Y}^* \mid (\mathbf{y}, \Lambda \mathbf{w}) + \langle \ell, \mathbf{w} \rangle = 0, \quad \forall \mathbf{w} \in \mathbf{V}_0\}$.

Note that since

$$(\mathbf{y}, \Lambda(\mathbf{u}_0 + \mathbf{w})) + \langle \ell, (\mathbf{u}_0 + \mathbf{w}) \rangle = (\mathbf{y}, \Lambda \mathbf{u}_0) + \langle \ell, \mathbf{u}_0 \rangle$$

we see that \mathbf{I}^* does not depend on the form of \mathbf{u}_0 inside Ω .

The problem dual to \mathcal{P} is as follows.

Problem \mathcal{P}^* . Find $\mathbf{p} \in \mathbf{Q}_\ell^*$ such that

$$\mathbf{I}^*(\mathbf{p}) = \sup_{\mathbf{y} \in \mathbf{Q}_\ell^*} \mathbf{I}^*(\mathbf{y}) := \sup \mathcal{P}^* \leq \inf \mathcal{P}.$$

The minimizer \mathbf{u} and the maximizer \mathbf{p} satisfy the conditions

$$\begin{aligned} (\mathcal{A}\Lambda\mathbf{u}, \Lambda\mathbf{w}) + \langle \ell, \mathbf{w} \rangle &= \mathbf{0} \quad \forall \mathbf{w} \in \mathbf{V}_0, \\ (\Lambda\mathbf{u}_0 - \mathcal{A}^{-1}\mathbf{p}, \mathbf{y}) &= \mathbf{0}, \quad \forall \mathbf{y} \in \mathbf{Q}_0^*, \end{aligned}$$

where $\mathbf{Q}_0^* := \{\mathbf{y} \in \mathbf{Y}^* \mid (\mathbf{y}, \Lambda\mathbf{w}) = \mathbf{0}, \quad \forall \mathbf{w} \in \mathbf{V}_0\}$.

We see that $\mathcal{A}\Lambda\mathbf{u} \in \mathbf{Q}_\ell^*$.

Take

$$I^*(\mathcal{A}\mathbf{\Lambda}\mathbf{u}) = (\mathcal{A}\mathbf{\Lambda}\mathbf{u}, \mathbf{\Lambda}\mathbf{u}_0) - \frac{1}{2} \|\mathcal{A}\mathbf{\Lambda}\mathbf{u}\|^2 + \langle \ell, \mathbf{u}_0 \rangle.$$

Recall that the dual functional does not depend on \mathbf{u}_0 inside Ω . Therefore, we set $\mathbf{u}_0 = \mathbf{u}$ and observe that

$$I^*(\mathcal{A}\mathbf{\Lambda}\mathbf{u}) = (\mathcal{A}\mathbf{\Lambda}\mathbf{u}, \mathbf{\Lambda}\mathbf{u}) - \frac{1}{2} \|\mathcal{A}\mathbf{\Lambda}\mathbf{u}\|_*^2 + \langle \ell, \mathbf{u} \rangle \leq \sup \mathcal{P}^*.$$

Since $\|\mathcal{A}\mathbf{\Lambda}\mathbf{u}\|_*^2 = (\mathcal{A}^{-1}\mathcal{A}\mathbf{\Lambda}\mathbf{u}, \mathcal{A}\mathbf{\Lambda}\mathbf{u}) = \|\mathbf{\Lambda}\mathbf{u}\|^2$, we see that

$$I^*(\mathcal{A}\mathbf{\Lambda}\mathbf{u}) = \mathbf{J}(\mathbf{u}) = \inf \mathcal{P}$$

Thus

$$\sup \mathcal{P}^* = \inf \mathcal{P}$$

The relation $\mathbf{l}^*(\mathbf{p}) = \mathbf{J}(\mathbf{u})$ means that

$$(\mathbf{p}, \Lambda \mathbf{u}) - \frac{1}{2} \|\mathbf{p}\|_*^2 + \langle \ell, \mathbf{u} \rangle = \frac{1}{2} \|\Lambda \mathbf{u}\|^2 + \langle \ell, \mathbf{u} \rangle,$$

which is equivalent to the relation

$$\mathbf{D}(\Lambda \mathbf{u}, \mathbf{p}) = \frac{1}{2} \|\Lambda \mathbf{u}\|^2 + \frac{1}{2} \|\mathbf{p}\|_*^2 - (\mathbf{p}, \Lambda \mathbf{u}) = 0.$$

From the above we see that $\Lambda \mathbf{u}$ and \mathbf{p} are joined by a certain relation:

$$\mathbf{p} = \mathcal{A} \Lambda \mathbf{u}$$

This is the so-called **duality relation** for the pair (\mathbf{u}, \mathbf{p}) .

Let $v \in V_0 + u_0$ and $\mathbf{y} \in \mathbf{Y}^*$ be some approximations of \mathbf{u} and \mathbf{p} , respectively. Our goal is to obtain two-sided estimates of the quantities $\|\Lambda(\mathbf{v} - \mathbf{u})\|$ and $\|\mathbf{y} - \mathbf{p}\|_*$ that are norms of **deviations** from the exact solutions \mathbf{u} and \mathbf{p} . First, we establish the following basic result.

Theorem

For any $v \in \mathbf{V}_0 + \mathbf{u}_0$ and $\mathbf{q} \in \mathbf{Q}_\ell^*$,

$$\|\Lambda(\mathbf{v} - \mathbf{u})\|^2 + \|\mathbf{q} - \mathbf{p}\|_*^2 = 2(\mathbf{J}(\mathbf{v}) - \mathbf{I}^*(\mathbf{q})), \quad (74)$$

$$\|\Lambda(\mathbf{v} - \mathbf{u})\|^2 + \|\mathbf{q} - \mathbf{p}\|_*^2 = 2\mathbf{D}(\Lambda\mathbf{v}, \mathbf{q}). \quad (75)$$

Proof

By the stationarity relations, we have

$$\begin{aligned} \frac{1}{2} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2 &= \frac{1}{2} (\mathcal{A}\Lambda(\mathbf{v} - \mathbf{u}), \Lambda(\mathbf{v} - \mathbf{u})) = \mathbf{J}(\mathbf{v}) - \mathbf{J}(\mathbf{u}) + \\ &+ (\mathcal{A}\Lambda\mathbf{u}, \Lambda(\mathbf{u} - \mathbf{v})) + \langle \ell, \mathbf{u} - \mathbf{v} \rangle = \\ &= \mathbf{J}(\mathbf{v}) - \mathbf{J}(\mathbf{u}). \end{aligned}$$

Analogously

$$\begin{aligned} \frac{1}{2} \|\mathbf{q} - \mathbf{p}\|_*^2 &= \frac{1}{2} (\mathcal{A}^{-1}(\mathbf{p} - \mathbf{q}), \mathbf{p} - \mathbf{q}) = \\ &= \mathbf{I}^*(\mathbf{p}) - \mathbf{I}^*(\mathbf{q}) - (\Lambda\mathbf{u}_0 - \mathcal{A}^{-1}\mathbf{p}, \mathbf{p} - \mathbf{q}) = \mathbf{I}^*(\mathbf{p}) - \mathbf{I}^*(\mathbf{q}). \end{aligned}$$

Since $\mathbf{J}(\mathbf{u}) = \mathbf{I}^*(\mathbf{p})$, we sum two relations and obtain (74). For $\mathbf{q} \in \mathbf{Q}_\ell^*$ the difference $\mathbf{J}(\mathbf{v}) - \mathbf{I}^*(\mathbf{q})$ is equal to $\mathbf{D}(\Lambda\mathbf{v}, \mathbf{q})$, so that (75) follows from (74).

The estimates (74) and (75) are valid only for $\mathbf{q} \in \mathbf{Q}_\ell^*$, which poses some difficulties. Below it is shown how we can overcome this drawback. First, we establish one subsidiary result.

Theorem

Let $\mathbf{q} \in \mathbf{Q}_\ell^*$, $\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0$, $\beta \in \mathbb{R}_+$, and $\mathbf{y} \in \mathbf{Y}^*$. Then

$$\mathbf{J}(\mathbf{v}) - \mathbf{I}^*(\mathbf{q}) \leq (\mathbf{1} + \beta)\mathbf{D}(\mathbf{\Lambda}\mathbf{v}, \mathbf{y}) + \frac{\mathbf{1} + \beta}{2\beta} \|\mathbf{q} - \mathbf{y}\|_*^2. \quad (76)$$

Note that

$$\begin{aligned} \mathbf{D}(\mathbf{\Lambda}\mathbf{v}, \mathbf{y}) &= \frac{\mathbf{1}}{2}(\mathcal{A}\mathbf{\Lambda}\mathbf{v}, \mathbf{\Lambda}\mathbf{v}) + \frac{\mathbf{1}}{2}(\mathcal{A}^{-1}\mathbf{p}, \mathbf{p}) - (\mathbf{y}, \mathbf{\Lambda}\mathbf{u}) = \\ &= (\mathcal{A}\mathbf{\Lambda}\mathbf{v} - \mathbf{y}, \mathbf{\Lambda}\mathbf{v} - \mathcal{A}^{-1}\mathbf{y}) = \\ &= (\mathcal{A}(\mathbf{\Lambda}\mathbf{v} - \mathcal{A}^{-1}\mathbf{y}), \mathbf{\Lambda}\mathbf{v} - \mathcal{A}^{-1}\mathbf{y}) = \|\mathbf{\Lambda}\mathbf{v} - \mathcal{A}^{-1}\mathbf{y}\|. \end{aligned}$$

Proof

For any $\mathbf{y} \in \mathbf{Y}^*$, we have

$$\begin{aligned} \mathbf{J}(\mathbf{v}) - \mathbf{I}^*(\mathbf{q}) &= \frac{1}{2} \left(\|\Lambda \mathbf{v}\|^2 + \|\mathbf{y}\|_*^2 \right) + \\ &\quad + \frac{1}{2} \left(\|\mathbf{q}\|_*^2 - \|\mathbf{y}\|_*^2 \right) - (\Lambda \mathbf{u}_0, \mathbf{q}) + \langle \ell, \mathbf{v} - \mathbf{u}_0 \rangle. \end{aligned}$$

Since $\langle \ell, \mathbf{v} - \mathbf{u}_0 \rangle = (\mathbf{q}, \Lambda(\mathbf{u}_0 - \mathbf{v}))$, we find that

$$\begin{aligned} \mathbf{J}(\mathbf{v}) - \mathbf{I}^*(\mathbf{q}) &= \frac{1}{2} \left(\|\Lambda \mathbf{v}\|^2 + \|\mathbf{y}\|_*^2 \right) + \frac{1}{2} \left(\|\mathbf{q}\|_*^2 - \|\mathbf{y}\|_*^2 \right) - (\mathbf{q}, \Lambda \mathbf{v}) = \\ &= \mathbf{D}(\Lambda \mathbf{v}, \mathbf{y}) + (\mathbf{y} - \mathbf{q}, \Lambda \mathbf{v} - \mathcal{A}^{-1} \mathbf{y}) + \frac{1}{2} \|\mathbf{q} - \mathbf{y}\|_*^2. \end{aligned}$$

This relation yields (76) if we use the Young's inequality

$$2(\mathbf{y} - \mathbf{q}, \Lambda \mathbf{v} - \mathcal{A}^{-1} \mathbf{y}) \leq \beta \|\Lambda \mathbf{v} - \mathcal{A}^{-1} \mathbf{y}\|^2 + \beta^{-1} \|\mathbf{y} - \mathbf{q}\|_*^2.$$

Another form of the estimate

Introduce the quantity

$$\mathbf{d}_\ell^2(\mathbf{y}) := \inf_{\mathbf{q} \in \mathbf{Q}_\ell^*} \|\mathbf{q} - \mathbf{y}\|_*^2,$$

which is the distance to \mathbf{Q}_ℓ^* . Then, (76) imply the estimate

$$\frac{1}{2} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2 \leq (1 + \beta) \mathbf{D}(\Lambda \mathbf{v}, \mathbf{y}) + \left(1 + \frac{1}{\beta}\right) \frac{1}{2} \mathbf{d}_\ell^2(\mathbf{y})$$

where $\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0$ and $\mathbf{y} \in \mathbf{Y}^*$. We rewrite this estimate as

$$\frac{1}{2} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2 \leq \mathcal{M}(\mathbf{v}, \beta), \quad \forall \mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0, \quad \beta \in \mathbb{R}_+,$$

where

$$\mathcal{M}(\mathbf{v}, \beta) := \inf_{\mathbf{y} \in \mathbf{Y}^*} \left\{ (1 + \beta) \mathbf{D}(\Lambda \mathbf{v}, \mathbf{y}) + \left(1 + \frac{1}{\beta}\right) \frac{1}{2} \mathbf{d}_\ell^2(\mathbf{y}) \right\}.$$

Now, we proceed to finding computable upper bounds for the quantity \mathbf{d}_ℓ .
The first step is given by

Theorem

$$\frac{1}{2} \mathbf{d}_\ell^2(\mathbf{y}) = \sup_{\mathbf{w} \in \mathbf{V}_0} \left\{ -\frac{1}{2} \|\boldsymbol{\Lambda} \mathbf{w}\|^2 - \langle \ell, \mathbf{w} \rangle - (\mathbf{y}, \boldsymbol{\Lambda} \mathbf{w}) \right\}.$$

Proof

This assertion comes from that $\inf \mathcal{P} = \sup \mathcal{P}^*$. Indeed,

$$\frac{1}{2} \mathbf{d}_\ell^2(\mathbf{y}) = - \sup_{\boldsymbol{\eta}^* \in \mathbf{Q}_\ell^*} \left\{ -\frac{1}{2} \|\mathbf{y} - \boldsymbol{\eta}^*\|_*^2 \right\} = - \sup_{\boldsymbol{\eta}^* \in \mathbf{Q}_\ell^* - \mathbf{y}} \left\{ -\frac{1}{2} \|\boldsymbol{\eta}^*\|_*^2 \right\},$$

where $\mathbf{Q}_\ell^* - \mathbf{y} := \{\boldsymbol{\eta}^* \in \mathbf{Y}^* \mid \boldsymbol{\eta}^* = \boldsymbol{\alpha}^* - \mathbf{y}, \quad \boldsymbol{\alpha}^* \in \mathbf{Q}_\ell^*\}$.
In other words, $\boldsymbol{\eta}^* \in \mathbf{Q}_\ell^* - \mathbf{y}$ if

$$(\boldsymbol{\eta}^*, \boldsymbol{\Lambda} \mathbf{w}) = -\langle \boldsymbol{\ell}, \mathbf{w} \rangle - (\mathbf{y}, \boldsymbol{\Lambda} \mathbf{w}), \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

The right-hand side of this relation is a linear continuous functional. We denote it by $\boldsymbol{\ell}_y$ and rewrite the relation as follows:

$$(\boldsymbol{\eta}^*, \boldsymbol{\Lambda} \mathbf{w}) + \langle \boldsymbol{\ell}_y, \mathbf{w} \rangle = \mathbf{0} \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

Then, $\mathbf{Q}_\ell^* - \mathbf{y} = \mathbf{Q}_{\ell_y}^*$ and

$$\frac{1}{2} \mathbf{d}_\ell^2(\mathbf{y}) = - \sup_{\boldsymbol{\eta}^* \in \mathbf{Q}_{\ell_y}^*} \left\{ -\frac{1}{2} \|\boldsymbol{\eta}^*\|_*^2 \right\}.$$

This maximization problem is a form of **Problem \mathcal{P}^*** if set $\mathbf{u}_0 = \mathbf{0}$ and $\ell = \ell_y$. Since $\sup \mathcal{P}^* = \inf \mathcal{P}$, we have

$$\begin{aligned} \frac{1}{2} \mathbf{d}_\ell^2(\mathbf{y}) &= - \inf_{\mathbf{w} \in \mathbf{V}_0} \left\{ \frac{1}{2} \|\boldsymbol{\Lambda} \mathbf{w}\|^2 + \langle \ell_y, \mathbf{w} \rangle \right\} = \\ &= - \inf_{\mathbf{w} \in \mathbf{V}_0} \left\{ \frac{1}{2} \|\boldsymbol{\Lambda} \mathbf{w}\|^2 + \langle \ell, \mathbf{w} \rangle + (\mathbf{y}, \boldsymbol{\Lambda} \mathbf{w}) \right\} = \\ &= \sup_{\mathbf{w} \in \mathbf{V}_0} \left\{ -\frac{1}{2} \|\boldsymbol{\Lambda} \mathbf{w}\|^2 - \langle \ell, \mathbf{w} \rangle - (\mathbf{y}, \boldsymbol{\Lambda} \mathbf{w}) \right\}. \end{aligned}$$

□

Corollary

We arrive at the conclusion that the majorant $\mathcal{M}(\mathbf{v}, \beta)$ has a minimax form

$$\mathcal{M}(\mathbf{v}, \beta) = \inf_{\mathbf{y} \in \mathbf{Y}^*} \sup_{\mathbf{w} \in \mathbf{V}_0} \left\{ (1 + \beta) \mathbf{D}(\boldsymbol{\Lambda} \mathbf{v}, \mathbf{y}) + \frac{1 + \beta}{\beta} \left(-(\mathbf{y}, \boldsymbol{\Lambda} \mathbf{w}) - \mathbf{J}(\mathbf{w}) \right) \right\}. \quad (77)$$

Further, we use (77) for deriving upper and lower error bounds.

Upper estimates of $\| \mathbf{v} - \mathbf{u} \|$

In the relation

$$\mathcal{M}(\mathbf{v}, \beta) \leq (\mathbf{1} + \beta) \mathbf{D}(\Lambda \mathbf{v}, \mathbf{y}) + \left(1 + \frac{1}{\beta}\right) \sup_{\mathbf{w} \in \mathbf{V}_0} \left\{ -\frac{1}{2} \| \Lambda \mathbf{w} \|^2 - \langle \ell, \mathbf{w} \rangle - (\mathbf{y}, \Lambda \mathbf{w}) \right\},$$

we will estimate the value of supremum. Since Λ^* is the operator conjugate to Λ , i.e.,

$$(\mathbf{y}, \Lambda \mathbf{w}) = \langle \Lambda^* \mathbf{y}, \mathbf{w} \rangle, \quad \forall \mathbf{w} \in \mathbf{V}_0,$$

we have

$$\langle \ell, \mathbf{w} \rangle + (\mathbf{y}, \Lambda \mathbf{w}) = \langle \ell + \Lambda^* \mathbf{y}, \mathbf{w} \rangle \leq \| \ell + \Lambda^* \mathbf{y} \| \| \Lambda \mathbf{w} \|.$$

Here

$$\|\ell + \Lambda^* \mathbf{y}\| := \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\langle \ell + \Lambda^* \mathbf{y}, \mathbf{w} \rangle}{\|\Lambda \mathbf{w}\|} < +\infty.$$

We see that

$$\begin{aligned} & \sup_{\mathbf{w} \in \mathbf{V}_0} \left\{ -\frac{1}{2} \|\Lambda \mathbf{w}\|^2 - \langle \ell, \mathbf{w} \rangle - (\mathbf{y}, \Lambda \mathbf{w}) \right\} \leq \\ \leq & \sup_{\mathbf{w} \in \mathbf{V}_0} \left\{ -\frac{1}{2} \|\Lambda \mathbf{w}\|^2 + \|\ell + \Lambda^* \mathbf{y}\| \|\Lambda \mathbf{w}\| \right\} \leq \\ \leq & \sup_{t > 0} \left\{ -\frac{1}{2} t^2 + \|\ell + \Lambda^* \mathbf{y}\| t \right\} = \frac{1}{2} \|\ell + \Lambda^* \mathbf{y}\|^2. \end{aligned}$$

Thus, we obtain

$$\frac{1}{2} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2 \leq (1 + \beta) \mathbf{D}(\Lambda \mathbf{v}, \mathbf{y}) + \frac{1 + \beta}{2\beta} \|\ell + \Lambda^* \mathbf{y}\|^2. \quad (78)$$

Deviation Majorant for the problem $\Lambda^* \mathcal{A} \Lambda u + \ell = 0$

$$\begin{aligned}
 (\mathcal{A} \Lambda(\mathbf{v} - \mathbf{u}), \Lambda(\mathbf{v} - \mathbf{u})) &\leq \\
 &\leq (1 + \beta) \left((\mathcal{A} \Lambda \mathbf{v}, \Lambda \mathbf{v}) + (\mathcal{A}^{-1} \mathbf{y}, \mathbf{y}) - 2(\mathbf{y}, \Lambda \mathbf{v}) \right) + \\
 &\quad + \frac{1 + \beta}{\beta} \mathbf{c}^2 \|\ell + \Lambda^* \mathbf{y}\|^2.
 \end{aligned}$$

In the above, $\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0$, $\beta > 0$, $\mathbf{y} \in \mathbf{U}$.

Theorem

For any $v \in \mathbf{V}_0 + \mathbf{u}_0$,

$$\frac{1}{2} \|\Lambda(\mathbf{u} - \mathbf{v})\|^2 = \inf_{\substack{\mathbf{y} \in \mathbf{Q}^* \\ \beta > 0}} \mathbf{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y}).$$

If $\ell \in \mathbf{U}$, then $\mathbf{p} \in \mathbf{Q}^*$ and, therefore,

$$\inf_{\substack{\mathbf{y} \in \mathbf{Q}^* \\ \beta > 0}} \mathbf{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y}) \leq \mathbf{M}_{\oplus}(\mathbf{v}, \varepsilon, \mathbf{p}) = (1 + \varepsilon) \frac{1}{2} \|\Lambda(\mathbf{u} - \mathbf{v})\|^2,$$

where $\varepsilon > 0$ may be taken arbitrarily small.

Hence, the majorant \mathbf{M}_{\oplus} is **reliable** and **exact**.

Lower estimates

Recall the minimax form of the Majorant

$$\mathcal{M}(\mathbf{v}, \beta) = \inf_{\mathbf{y} \in \mathbf{Y}^*} \sup_{\mathbf{w} \in \mathbf{V}_0} \left\{ (1 + \beta) \mathbf{D}(\Lambda \mathbf{v}, \mathbf{y}) + \frac{1 + \beta}{\beta} \left(-(\mathbf{y}, \Lambda \mathbf{w}) - \mathbf{J}(\mathbf{w}) \right) \right\}.$$

Since $\sup \inf \leq \inf \sup$, we have

$$\begin{aligned} \mathcal{M}(\mathbf{v}, \beta) \geq \sup_{\mathbf{w} \in \mathbf{V}_0} \inf_{\mathbf{y} \in \mathbf{Y}^*} \left\{ (1 + \beta) \mathbf{D}(\Lambda \mathbf{v}, \mathbf{y}) - \right. \\ \left. - \left(1 + \frac{1}{\beta} \right) \left(\frac{1}{2} \|\Lambda \mathbf{w}\|^2 + \langle \ell, \mathbf{w} \rangle + (\mathbf{y}, \Lambda \mathbf{w}) \right) \right\}. \end{aligned}$$

Thus, for any $\mathbf{w} \in \mathbf{V}_0$

$$\begin{aligned} \mathcal{M}(\mathbf{v}, \beta) \geq & \inf_{\mathbf{y} \in \mathbf{Y}^*} \left\{ (1 + \beta) \left(\frac{1}{2} \|\mathbf{y}\|_*^2 - (\mathbf{y}, \Lambda \mathbf{v}) \right) - \left(1 + \frac{1}{\beta} \right) (\mathbf{y}, \Lambda \mathbf{w}) \right\} + \\ & + (1 + \beta) \frac{1}{2} \|\Lambda \mathbf{v}\|^2 - \left(1 + \frac{1}{\beta} \right) \left(\frac{1}{2} \|\Lambda \mathbf{w}\|^2 + \langle \ell, \mathbf{w} \rangle \right), \end{aligned}$$

Evidently, this estimate is also valid for the function $\beta \mathbf{w}$, which yields

$$\begin{aligned} \mathcal{M}(\mathbf{v}, \beta) \geq (1 + \beta) \inf_{\mathbf{y} \in \mathbf{Y}^*} \left\{ \frac{1}{2} \|\mathbf{y}\|_*^2 - (\mathbf{y}, \Lambda(\mathbf{v} + \mathbf{w})) \right\} + \\ + (1 + \beta) \left(\frac{1}{2} \|\Lambda \mathbf{v}\|^2 - \frac{\beta}{2} \|\Lambda \mathbf{w}\|^2 - \langle \ell, \mathbf{w} \rangle \right). \end{aligned}$$

Note that

$$\begin{aligned} \inf_{\mathbf{y} \in \mathbf{Y}^*} \left\{ \frac{1}{2} \|\mathbf{y}\|_*^2 - (\mathbf{y}, \Lambda(\mathbf{v} + \mathbf{w})) \right\} &\geq \\ &\geq \inf_{\mathbf{y} \in \mathbf{Y}^*} \left\{ \frac{1}{2} \|\mathbf{y}\|_*^2 - \|\mathbf{y}\|_* \|\Lambda(\mathbf{v} + \mathbf{w})\| \right\} = -\frac{1}{2} \|\Lambda(\mathbf{v} + \mathbf{w})\|^2. \end{aligned}$$

Thus, we obtain

$$\begin{aligned} \mathcal{M}(\mathbf{v}, \beta) &\geq (1 + \beta) \left\{ -\frac{1}{2} \|\Lambda(\mathbf{v} + \mathbf{w})\|^2 + \right. \\ &\quad \left. + \frac{1}{2} \|\Lambda\mathbf{v}\|^2 - \frac{\beta}{2} \|\Lambda\mathbf{w}\|^2 - \langle \ell, \mathbf{w} \rangle \right\} = \\ &= (1 + \beta) \left\{ -(\mathcal{A}\Lambda\mathbf{v}, \Lambda\mathbf{w}) - \frac{1 + \beta}{2} \|\Lambda\mathbf{w}\|^2 - \langle \ell, \mathbf{w} \rangle \right\}. \end{aligned}$$

In

$$(1 + \beta) \left\{ -(\mathcal{A}\Lambda\mathbf{v}, \Lambda\mathbf{w}) - \frac{1 + \beta}{2} \|\Lambda\mathbf{w}\|^2 - \langle \ell, \mathbf{w} \rangle \right\}.$$

\mathbf{w} is an arbitrary function in \mathbf{V}_0 . We may replace

$$\mathbf{w} \quad \text{by} \quad \frac{\mathbf{w}}{1 + \beta}.$$

Such a replacement leads to the **Minorant** $M_{\ominus}(\mathbf{v}, \mathbf{w})$ that gives a **lower bound** of the deviation from exact solution:

For any $\mathbf{w} \in \mathbf{V}_0$,

$$\frac{1}{2} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2 \geq -\frac{1}{2} \|\Lambda\mathbf{w}\|^2 - (\mathcal{A}\Lambda\mathbf{v}, \Lambda\mathbf{w}) - \langle \ell, \mathbf{w} \rangle \quad (79)$$

Minorant is sharp

It is easy to see that

$$\sup_{\mathbf{w} \in \mathbf{V}_0} \mathbf{M}_\ominus(\mathbf{v}, \mathbf{w}) = \frac{1}{2} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2.$$

Indeed, take $\mathbf{w} = \mathbf{u} - \mathbf{v}$.

$$\mathbf{M}_\ominus(\mathbf{v}, \mathbf{u} - \mathbf{v}) = -\frac{1}{2} \|\Lambda(\mathbf{u} - \mathbf{v})\|^2 - (\mathcal{A}\Lambda\mathbf{v}, \Lambda\mathbf{u} - \mathbf{v}) - \langle \ell, \mathbf{u} - \mathbf{v} \rangle.$$

Represent the last two terms as follows:

$$\begin{aligned} & -(\mathcal{A}\Lambda\mathbf{v}, \Lambda(\mathbf{u} - \mathbf{v})) - \langle \ell, \mathbf{u} - \mathbf{v} \rangle = \\ & = -(\mathcal{A}\Lambda\mathbf{v}, \Lambda(\mathbf{u} - \mathbf{v})) + (\mathcal{A}\Lambda\mathbf{u}, \Lambda(\mathbf{u} - \mathbf{v})) = \\ & = (\mathcal{A}\Lambda(\mathbf{u} - \mathbf{v}), \Lambda(\mathbf{u} - \mathbf{v})) = \|\Lambda(\mathbf{u} - \mathbf{v})\|^2 \end{aligned}$$

so that this choice of \mathbf{w} gives the true error.

Remark.

We outline that for the exact solution $\mathbf{M}_\ominus = \mathbf{M}_\oplus = \mathbf{0}$! Indeed, assume that \mathbf{v} coincides with \mathbf{u} . In this case,

$$\mathbf{M}_\ominus(\mathbf{u}, \mathbf{w}) = -\frac{1}{2} \|\boldsymbol{\Lambda} \mathbf{w}\|^2 - (\mathcal{A} \boldsymbol{\Lambda} \mathbf{u}, \boldsymbol{\Lambda} \mathbf{w}) - \langle \boldsymbol{\ell}, \mathbf{w} \rangle = -\frac{1}{2} \|\boldsymbol{\Lambda} \mathbf{w}\|^2$$

and, therefore,

$$\sup_{\mathbf{w} \in \mathbf{V}_0} \mathbf{M}_\ominus(\mathbf{u}, \mathbf{w}) = \mathbf{0}.$$

The same is true for the majorant. Indeed, set $\hat{\mathbf{y}} = \mathcal{A} \boldsymbol{\Lambda} \mathbf{u}$. Then,

$$\mathbf{M}_\oplus(\mathbf{u}, \beta, \hat{\mathbf{y}}) = (1 + \beta) \mathbf{D}(\boldsymbol{\Lambda} \mathbf{u}, \hat{\mathbf{y}}) + \frac{1 + \beta}{2\beta} \mathbf{c}^2 \|\boldsymbol{\ell} + \boldsymbol{\Lambda}^* \mathcal{A} \boldsymbol{\Lambda} \mathbf{u}\|^2 = \mathbf{0}.$$

Thus,

$$\inf_{\mathbf{y} \in \mathbf{Y}^*} \mathbf{M}_\oplus(\mathbf{u}, \beta, \mathbf{y}) = \mathbf{0}.$$

Computability of two-sided estimates

By **computability** we mean that upper and lower estimates can be computed with any a priori given accuracy by solving finite-dimensional problems. In the case considered, they are certain problems for quadratic type integral functionals whose minimization (maximization) is performed by well-known methods.

Let $\{\mathbf{Y}_i^*\}_{i=1}^\infty$ and $\{\mathbf{V}_{0i}\}_{i=1}^\infty$ be two sequences of finite-dimensional subspaces that are dense in \mathbf{Q}^* and \mathbf{V}_0 , respectively, i.e., for any given $\varepsilon > 0$ and arbitrary elements $\mathbf{y} \in \mathbf{Y}^*$ and $\mathbf{w} \in \mathbf{V}_0$, one can find a natural number \mathbf{k}_ε such that

$$\inf_{\tilde{\mathbf{w}} \in \mathbf{V}_{0i}} \|\tilde{\mathbf{w}} - \mathbf{w}\|_{\mathbf{V}} \leq \varepsilon, \quad \inf_{\tilde{\mathbf{y}} \in \mathbf{Y}_i^*} \|\tilde{\mathbf{y}} - \mathbf{y}\|_{\mathbf{Q}^*} \leq \varepsilon, \quad \forall i \geq \mathbf{k}_\varepsilon.$$

Let us show that sequences of two-sided bounds converging to the actual error can be evaluated by minimizing the Majorant on $\{\mathbf{Y}_i^*\}$ and maximizing the Minorant on $\{\mathbf{V}_{0i}\}$.

Take a small $\varepsilon > 0$,. Then there exists a number \mathbf{k} and elements $\mathbf{w}_k \in \mathbf{V}_{0k}$ and $\mathbf{p}_k \in \mathbf{Y}_{0k}^*$ satisfying the conditions

$$\|\mathbf{w}_k - (\mathbf{u} - \mathbf{v})\|_{\mathbf{V}} \leq \varepsilon, \quad \|\mathbf{p}_k - \mathbf{p}\|_{\mathbf{Q}^*} \leq \varepsilon.$$

Define two quantities defined by solving **finite-dimensional problems**, namely

$$\mathbf{M}_{\oplus}^k = \inf_{\substack{\mathbf{y}_k \in \mathbf{Y}_k^* \\ \beta \in \mathbb{R}_+}} \mathbf{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y}_k), \quad \mathbf{M}_{\ominus}^k = \sup_{\mathbf{w}_k \in \mathbf{V}_{0k}} \mathbf{M}_{\ominus}(\mathbf{v}, \mathbf{w}_k).$$

By the definition

$$\mathbf{M}_{\ominus}(\mathbf{v}, \mathbf{w}_k) \leq \mathbf{M}_{\ominus}^k \leq \frac{1}{2} \|\mathbf{u} - \mathbf{v}\|^2 \leq \mathbf{M}_{\oplus}^k \leq \mathbf{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{p}_k).$$

The quantities M_{\ominus}^k and M_{\oplus}^k are **computable** (they require solving finite dimensional problems for quadratic type functionals). We will that

$$M_{\oplus}^k \rightarrow \frac{1}{2} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2,$$
$$M_{\ominus}^k \rightarrow \frac{1}{2} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2$$

as the dimensionality \mathbf{k} tends to $+\infty$.

Consider the **upper estimates**.

$$\begin{aligned} \mathbf{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{p}_k) &= (\mathbf{1} + \beta) \mathbf{D}(\Lambda \mathbf{v}, \mathbf{p}_k) + \frac{\mathbf{1} + \beta}{2\beta} \mathbf{c}^2 \|\ell + \Lambda^* \mathbf{p}_k\|^2 = \\ &= (\mathbf{1} + \beta) \mathbf{D}(\Lambda \mathbf{v}, \mathbf{p}_k) + \frac{\mathbf{1} + \beta}{2\beta} \mathbf{c}^2 \|\Lambda^* (\mathbf{p}_k - \mathbf{p})\|^2. \end{aligned}$$

Here

$$\begin{aligned} \mathbf{D}(\Lambda \mathbf{v}, \mathbf{p}_k) &= \frac{\mathbf{1}}{2} (\Lambda \mathbf{v} - \mathcal{A}^{-1} \mathbf{p}_k, \mathcal{A} \Lambda \mathbf{v} - \mathbf{p}_k) = \\ &= \frac{\mathbf{1}}{2} \left(\Lambda (\mathbf{v} - \mathbf{u}) - \mathcal{A}^{-1} (\mathbf{p}_k - \mathbf{p}), \mathcal{A} \Lambda (\mathbf{v} - \mathbf{u}) - (\mathbf{p}_k - \mathbf{p}) \right) = \\ &= \frac{\mathbf{1}}{2} \left\| \Lambda (\mathbf{v} - \mathbf{u}) \right\|^2 + \left\| \mathbf{p}_k - \mathbf{p} \right\|_*^2 - (\Lambda (\mathbf{v} - \mathbf{u}), \mathbf{p}_k - \mathbf{p}). \end{aligned}$$

From the latter estimate we see that

$$\mathbf{D}(\Lambda \mathbf{v}, \mathbf{p}_k) \leq \frac{\mathbf{1}}{2} \left\| \Lambda (\mathbf{v} - \mathbf{u}) \right\|^2 + \varepsilon \left\| \Lambda (\mathbf{v} - \mathbf{u}) \right\| + \frac{\mathbf{1}}{2} \varepsilon^2. \quad (80)$$

Since

$$\|\Lambda^*(\mathbf{p}_k - \mathbf{p})\|_{\mathbf{Q}^*} \leq \varepsilon,$$

we find that

$$\begin{aligned} \mathbf{M}_{\oplus}^k &\leq \mathbf{M}_{\oplus}(\mathbf{v}, \varepsilon, \mathbf{p}_k) = \\ &= (1 + \varepsilon) \left(\frac{1}{2} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2 + \varepsilon \|\Lambda(\mathbf{v} - \mathbf{u})\| + \frac{1}{2} \varepsilon^2 \right) + \frac{1 + \varepsilon}{2\varepsilon} \mathbf{c}^2 \varepsilon^2 = \\ &= \frac{1}{2} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2 + \mathbf{c}_4 \varepsilon + \mathbf{o}(\varepsilon^2). \end{aligned}$$

where $\mathbf{c}_4 = \frac{1}{2} (\mathbf{c} + 2 \|\Lambda(\mathbf{v} - \mathbf{u})\| + \|\Lambda(\mathbf{v} - \mathbf{u})\|^2)$. Thus, we conclude that

$$\mathbf{M}_{\oplus}^k \longrightarrow \frac{1}{2} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2 \quad \text{as } k \rightarrow \infty.$$

Remark

It is worth noting that the constant c_4 in the convergence term with ε depends on the norm of $(\mathbf{v} - \mathbf{u})$, so that we can await that for a good approximation convergence of the upper bounds to the exact value of the error is faster than in the case where $\|\mathbf{v} - \mathbf{u}\|$ is considerable. This phenomenon was observed in many numerical experiments. In general, finding an upper bound for a precise approximation takes less CPU time than for a coarse one.

Consider the **lower estimates**.

$$\begin{aligned}
 \mathbf{M}_{\ominus}(\mathbf{v}, \mathbf{w}_k) &= -\frac{1}{2} \|\Lambda \mathbf{w}_k\|^2 - (\mathcal{A}\Lambda \mathbf{v}, \Lambda \mathbf{w}_k) - \langle \ell, \mathbf{w}_k \rangle = \\
 &= -\frac{1}{2} \|\Lambda \mathbf{w}_k\|^2 + (\mathcal{A}\Lambda(\mathbf{u} - \mathbf{v}), \Lambda \mathbf{w}_k) = \\
 &= \frac{1}{2} \|\Lambda(\mathbf{u} - \mathbf{v})\|^2 - \frac{1}{2} \|\Lambda(\mathbf{w}_k - (\mathbf{u} - \mathbf{v}))\|^2 \geq \\
 &\geq \frac{1}{2} \|\Lambda(\mathbf{u} - \mathbf{v})\|^2 - \frac{1}{2} \mathbf{c}_2 \|\Lambda(\mathbf{w}_k - (\mathbf{u} - \mathbf{v}))\|^2.
 \end{aligned}$$

This implies the estimate

$$\frac{1}{2} \|\Lambda(\mathbf{u} - \mathbf{v})\|^2 \geq \mathbf{M}_{\ominus}^k \geq \frac{1}{2} \|\Lambda(\mathbf{u} - \mathbf{v})\|^2 - \mathbf{c}_5 \varepsilon^2,$$

where $\mathbf{c}_5 > \mathbf{0}$ depends on the norm of Λ . Thus,

$$\mathbf{M}_{\ominus}^k \rightarrow \frac{1}{2} \|\Lambda(\mathbf{u} - \mathbf{v})\|^2 \quad \text{as } k \rightarrow \infty.$$

Computable upper bound of the effectivity index

Having M_{\oplus}^k and M_{\ominus}^k , one can define the number

$$\eta_k := \frac{M_{\oplus}^k}{M_{\ominus}^k} \geq 1, \quad (81)$$

which gives an idea of the **quality of the error estimation**. From the above it follows that

$$\eta_k \rightarrow 1, \quad \text{as } k \rightarrow +\infty.$$

Relationships with other methods

$$\mathbf{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y}) = (1 + \beta)\mathbf{D}(\Lambda\mathbf{v}, \mathbf{y}) + \frac{1 + \beta}{2\beta} \mathbf{c}^2 \|\ell + \Lambda^*\mathbf{y}\|^2$$

involves an arbitrary function \mathbf{y} . We are aimed to show that some special choices of it lead to known error estimates.

We assume that $\langle \ell, \mathbf{w} \rangle = (\mathbf{g}, \mathbf{w})$, where $\mathbf{g} \in \mathbf{U}$, so that $\mathbf{p} \in \mathbf{Q}^* \subset \mathbf{Q}_{\ell}^*$ and

$$\mathbf{Q}_{\ell}^* := \{\mathbf{y} \in \mathbf{Q}^* \mid (\Lambda^*\mathbf{y} + \mathbf{g}, \mathbf{w}) = 0, \quad \forall \mathbf{w} \in \mathbf{V}_0\}.$$

Let us first define the function

$$\mathbf{y}_0 = \mathcal{A}\mathbf{v}. \quad (82)$$

A variety of options comes from the relation

$$\mathbf{y} = \mathbf{\Pi}\mathbf{y}_0, \quad (83)$$

where $\mathbf{\Pi}$ is a certain continuous mapping.

Residual based estimate

If Π is the identity mapping of \mathbf{Y}^* , i.e., $\mathbf{y} = \mathbf{y}_0 := \mathcal{A}\boldsymbol{\Lambda}\mathbf{v}$, then

$$\mathbf{D}(\boldsymbol{\Lambda}\mathbf{v}, \mathbf{y}_0^*) = \mathbf{0}.$$

Use the majorant in the form (78):

$$\frac{1}{2} \|\boldsymbol{\Lambda}(\mathbf{v} - \mathbf{u})\|^2 \leq (\mathbf{1} + \beta)\mathbf{D}(\boldsymbol{\Lambda}\mathbf{v}, \mathbf{y}) + \frac{\mathbf{1} + \beta}{2\beta} \|\boldsymbol{\ell} + \boldsymbol{\Lambda}^* \mathbf{y}\|^2.$$

Now, it contains only the second term, which after the minimization with respect to β gives

$$\|\boldsymbol{\Lambda}(\mathbf{v} - \mathbf{u})\| \leq \|\boldsymbol{\ell} + \boldsymbol{\Lambda}^* \mathcal{A}\boldsymbol{\Lambda}\mathbf{v}\| = \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{(\boldsymbol{\ell}, \mathbf{w}) + (\mathcal{A}\boldsymbol{\Lambda}\mathbf{v}, \boldsymbol{\Lambda}\mathbf{w})}{\|\boldsymbol{\Lambda}\mathbf{w}\|}. \quad (84)$$

If \mathbf{v} is obtained by FEM and $\mathbf{v} = \mathbf{u}_h \in \mathbf{V}_h := \mathbf{V}_{0h} + \mathbf{u}_0$, then we arrive at the following estimate:

$$\| \Lambda(\mathbf{u}_h - \mathbf{u}) \| \leq [\ell + \Lambda^* \mathcal{A} \Lambda \mathbf{u}_h] = \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{(\ell, \mathbf{w} - \mathbf{w}_h) + (\mathcal{A} \Lambda \mathbf{u}_h, \Lambda(\mathbf{w} - \mathbf{w}_h))}{\| \Lambda \mathbf{w} \|}.$$

We find an upper bound of the right-hand side by the arguments accepted in the classical residual method.

Conclusion

If in the functional a posteriori error estimate is applied to a FEM solution \mathbf{u}_h then we may select the variable \mathbf{y} in the simplest way as

$$\mathbf{y} = \Lambda \mathbf{u}_h.$$

Then, if \mathbf{u}_h is a Galerkin approximation, we can use this fact and obtain at an upper bound given by the *residual type a posteriori error estimate* that involves integral terms associated with finite elements and interelement jumps.

Estimates using post-processing of the dual variable

In $\mathbf{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y})$ the best choice is $\mathbf{y} = \mathbf{p} \in \mathbf{Q}^*$. Therefore, if $\mathbf{y}_0 \notin \mathbf{Q}^*$ then its mapping \mathbf{Q}^* could be a better approximation of \mathbf{p} . Let us denote such a mapping by Π_1 . We obtain

$$\mathbf{y}_1 = \Pi_1 \mathbf{y}_0 \in \mathbf{Q}^* \quad (85)$$

and the quantity $\mathbf{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y}_1)$, which leads to the error majorant

$$\begin{aligned} \mathbf{M}_{\oplus}^{(1)}(\mathbf{v}) = \inf_{\beta \in \mathbb{R}_+} \left\{ (1 + \beta) \mathbf{D}(\Lambda \mathbf{v}, \Pi_1(\mathcal{A} \Lambda \mathbf{v})) + \right. \\ \left. + \frac{1 + \beta}{2\beta} \mathbf{c}^2 \|\ell + \Lambda^* \Pi_1(\mathcal{A} \Lambda \mathbf{v})\|^2 \right\}. \quad (86) \end{aligned}$$

Particular case

In the simplest case associated with the problem

$$\Delta \mathbf{u} + \mathbf{f} = \mathbf{0}, \quad \mathbf{u} = \mathbf{u}_0 \quad \text{on } \partial\Omega$$

we have

$$\begin{aligned} \mathbf{M}_{\oplus}^{(1)}(\mathbf{u}_h) &= \\ &= \inf_{\beta \in \mathbb{R}_+} \left\{ (1+\beta) \|\nabla \mathbf{u}_h - \Pi_1(\nabla \mathbf{u}_h)\|^2 + \frac{(1+\beta) \mathbf{C}_{\Omega}^2}{2\beta} \|\mathbf{f} + \mathbf{div} \Pi_1(\nabla \mathbf{u}_h)\|^2 \right\}. \end{aligned}$$

If Π_1 is a gradient averaging operator, then the first term in the right-hand side is **the difference between the original and averaged gradient**, i.e. it coincides with a **gradient averaging indicator**. However, as we have seen in previous lectures, such an indicator cannot provide a reliable upper bound of the error. The second term in the right-hand side shows what is necessary to add **in order to provide the reliability**.

Estimates based on the "equilibration" of the dual variable

Let Π_2 maps \mathbf{Y}^* to the set \mathbf{Q}_ℓ^* . Define

$$\mathbf{y}_2 = \Pi_2 \mathbf{y}_0 \in \mathbf{Q}_\ell^*. \quad (87)$$

Then,

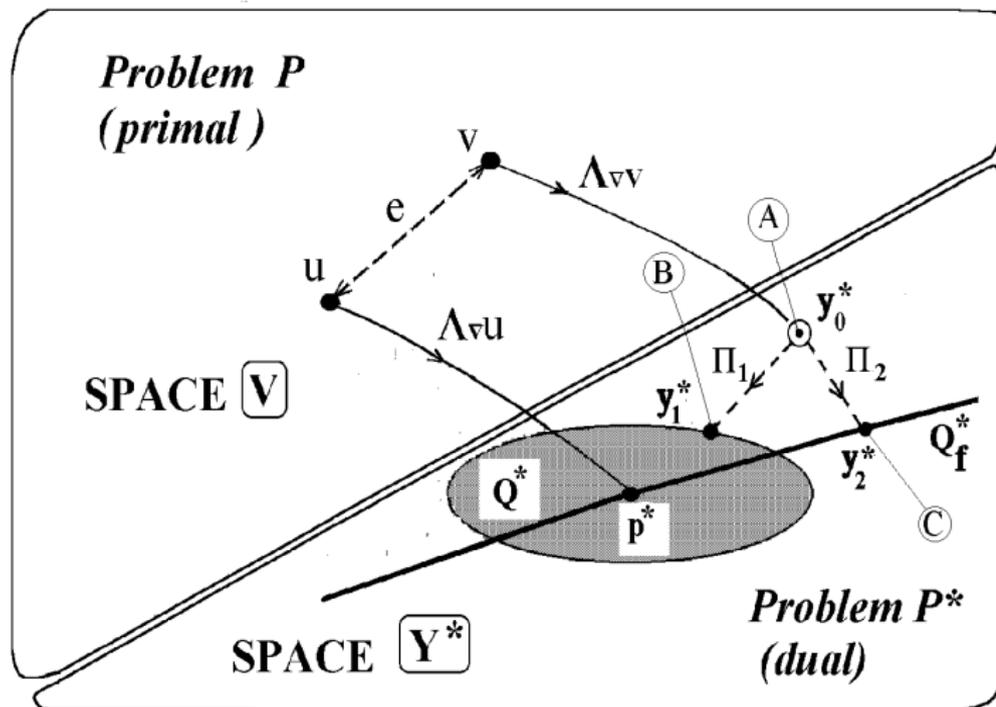
$$\Lambda^* \mathbf{y}_2 + \ell = \mathbf{0},$$

so that the Majorant has only the first term:

$$\mathbf{M}_{\oplus}^{(2)}(\mathbf{v}) = \mathbf{D}(\Lambda \mathbf{v}, \mathbf{y}_2).$$

Π_2 is natural to call an **equilibration operator**. In general, it is rather difficult to construct an "exact mapping" Π_2 to \mathbf{Q}_ℓ^* . One may use an operator $\tilde{\Pi}_2$, which provides an approximate "equilibration". In this case, the **second term of the Majorant does not vanish and should be taken into account**.

Various choices of the dual variable lead to certain a posteriori methods



A priori projection type error estimates

As an exercise, we now will derive classical a priori projection type error estimates from a functional a posteriori estimate. Let $\mathbf{u}_h \in \mathbf{V}_h$ be a Galerkin approximation of \mathbf{u} . We have

$$\|\Lambda(\mathbf{u} - \mathbf{u}_h)\|^2 \leq 2(1 + \beta) \mathbf{D}(\Lambda \mathbf{u}_h, \mathbf{y}) + \left(1 + \frac{1}{\beta}\right) \|\Lambda^* \mathbf{y} + \ell\|^2$$

Set here $\mathbf{y} = \mathcal{A}\Lambda \mathbf{v}_h$, where \mathbf{v}_h is an arbitrary element of \mathbf{V}_h . Then,

$$\begin{aligned} \|\Lambda^* \mathbf{y} + \ell\| &= \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{(\mathbf{y} - \mathbf{p}, \Lambda \mathbf{w})}{\|\Lambda \mathbf{w}\|} = \\ &= \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{(\mathcal{A}\Lambda(\mathbf{v}_h - \mathbf{u}), \Lambda \mathbf{w})}{\|\Lambda \mathbf{w}\|} \leq \|\Lambda(\mathbf{v}_h - \mathbf{u})\|. \end{aligned}$$

It is easy to see that

$$\mathbf{D}(\boldsymbol{\Lambda}\mathbf{u}_h, \mathcal{A}\boldsymbol{\Lambda}\mathbf{v}_h) = \mathbf{J}(\mathbf{v}_h) - \mathbf{J}(\mathbf{u}_h).$$

Indeed,

$$\begin{aligned} \mathbf{D}(\boldsymbol{\Lambda}\mathbf{u}_h, \mathcal{A}\boldsymbol{\Lambda}\mathbf{v}_h) &= \frac{1}{2}(\mathcal{A}\boldsymbol{\Lambda}\mathbf{v}_h, \boldsymbol{\Lambda}\mathbf{v}_h) + \langle \boldsymbol{\ell}, \mathbf{v}_h \rangle - \\ &\quad - \frac{1}{2}(\mathcal{A}\boldsymbol{\Lambda}\mathbf{u}_h, \boldsymbol{\Lambda}\mathbf{u}_h) - \langle \boldsymbol{\ell}, \mathbf{u}_h \rangle + \\ &\quad + (\mathcal{A}\boldsymbol{\Lambda}\mathbf{u}_h, \boldsymbol{\Lambda}(\mathbf{u}_h - \mathbf{v}_h)) + \langle \boldsymbol{\ell}, \mathbf{u}_h - \mathbf{v}_h \rangle. \end{aligned}$$

Since $\mathbf{u}_h \in \mathbf{V}_h$ is a Galerkin approximation, the last two terms vanish and we obtain the relation.

We know that

$$\begin{aligned} \|\boldsymbol{\Lambda}(\mathbf{u}_h - \mathbf{u})\|^2 &= 2(\mathbf{J}(\mathbf{u}_h) - \mathbf{J}(\mathbf{u})), \\ \|\boldsymbol{\Lambda}(\mathbf{v}_h - \mathbf{u})\|^2 &= 2(\mathbf{J}(\mathbf{v}_h) - \mathbf{J}(\mathbf{u})). \end{aligned}$$

Therefore,

$$\begin{aligned} 2\mathbf{D}(\mathbf{\Lambda}\mathbf{u}_h, \mathcal{A}\mathbf{\Lambda}\mathbf{v}_h) &= 2(\mathbf{J}(\mathbf{v}_h) - \mathbf{J}(\mathbf{u})) - 2(\mathbf{J}(\mathbf{u}_h) - \mathbf{J}(\mathbf{u})) = \\ &= \|\mathbf{\Lambda}(\mathbf{v}_h - \mathbf{u})\|^2 - \|\mathbf{\Lambda}(\mathbf{u}_h - \mathbf{u})\|^2. \end{aligned}$$

Now, the error estimate comes in the form

$$\begin{aligned} \|\mathbf{\Lambda}(\mathbf{u} - \mathbf{u}_h)\|^2 &\leq (1 + \beta)(\|\mathbf{\Lambda}(\mathbf{v}_h - \mathbf{u})\|^2 - \|\mathbf{\Lambda}(\mathbf{u}_h - \mathbf{u})\|^2) + \\ &+ \left(1 + \frac{1}{\beta}\right) \|\mathbf{\Lambda}(\mathbf{v}_h - \mathbf{u})\|^2. \end{aligned}$$

Thus, we obtain

$$\begin{aligned} (2 + \beta) \|\mathbf{\Lambda}(\mathbf{u} - \mathbf{u}_h)\|^2 &\leq \\ &\leq (1 + \beta) \|\mathbf{\Lambda}(\mathbf{v}_h - \mathbf{u})\|^2 + \left(1 + \frac{1}{\beta}\right) \|\mathbf{\Lambda}(\mathbf{v}_h - \mathbf{u})\|^2, \end{aligned}$$

We see that

$$\| \Lambda(\mathbf{u} - \mathbf{u}_h) \| \leq \left(1 + \frac{1}{\beta(2 + \beta)} \right) \| \Lambda(\mathbf{u} - \mathbf{v}_h) \| .$$

Since β is an arbitrary positive number, we arrive at the projection type error estimate

$$\| \Lambda(\mathbf{u} - \mathbf{u}_h) \| \leq \inf_{\mathbf{v}_h \in \mathbf{V}_h} \| \Lambda(\mathbf{u} - \mathbf{v}_h) \| .$$

APPLICATIONS TO PARTICULAR CLASSES OF PARTIAL DIFFERENTIAL EQUATIONS

Diffusion equation

Let \mathcal{A} is produced by a matrix $\mathbf{A} = \{\mathbf{a}_{ij}\} = \{\mathbf{a}_{ji}\}$, $\mathbf{V} = \mathbf{H}^1(\Omega)$, where Ω is a Lipschitz domain, $\mathbf{U} = \mathbf{L}^2(\Omega, \mathbb{R}^n)$, and $\mathbf{A}\mathbf{w} = \nabla\mathbf{w}$. Let the entries of \mathbf{A} be bounded at almost all points of Ω and such that

$$\mathbf{c}_1|\xi|^2 \leq \mathbf{a}_{ij}\xi_i\xi_j \leq \mathbf{c}_2|\xi|^2, \quad \forall \xi \in \mathbb{R}^n. \quad (88)$$

Then, the spaces \mathbf{Y} and \mathbf{Y}^* have the norms

$$\|\mathbf{y}\|^2 = \int_{\Omega} \mathbf{A}\mathbf{y} \cdot \mathbf{y} \, dx, \quad \|\mathbf{y}\|_*^2 = \int_{\Omega} \mathbf{A}^{-1}\mathbf{y} \cdot \mathbf{y} \, dx.$$

Dirichlet boundary conditions

We begin with the problem

$$\mathbf{div} \mathbf{A} \nabla \mathbf{u} = \mathbf{f} \quad \text{in} \quad \Omega, \quad (89)$$

$$\mathbf{u} = \mathbf{u}_0 \quad \text{on} \quad \partial\Omega. \quad (90)$$

In this case, $\mathbf{V}_0 = \mathring{\mathbf{H}}^1(\Omega)$ and \mathbf{u} meets the integral identity

$$\int_{\Omega} \mathbf{A} \nabla \mathbf{u} \cdot \nabla \mathbf{w} \, d\mathbf{x} + \langle \mathbf{f}, \mathbf{w} \rangle = 0, \quad \forall \mathbf{w} \in \mathbf{V}_0. \quad (91)$$

The relation $(\mathbf{y}, \mathbf{\Lambda} \mathbf{w}) = \langle \mathbf{\Lambda}^* \mathbf{y}, \mathbf{w} \rangle$ has the form

$$\int_{\Omega} \mathbf{y} \cdot \nabla \mathbf{w} \, d\mathbf{x} = \langle -\mathbf{div} \, \mathbf{y}, \mathbf{w} \rangle,$$

where $\mathbf{\Lambda}^* = -\mathbf{div}$ and $\mathbf{div} \, \mathbf{y}$ is in $\mathbf{H}^{-1}(\Omega)$.

The operator \mathbf{A} satisfies the required inequality

$$c_{\Omega} \|\nabla \mathbf{w}\| \geq \|\mathbf{w}\|, \quad \forall \mathbf{w} \in \mathring{\mathbf{H}}^1(\Omega).$$

Upper estimates of $\|\mathbf{v} - \mathbf{u}\|$ for an approximation $\mathbf{v} \in V_0 + u_0$ follow from the general estimate presented in Lecture 5. We have

$$\frac{1}{2} \int_{\Omega} \mathbf{A} \nabla(\mathbf{v} - \mathbf{u}) \cdot \nabla(\mathbf{v} - \mathbf{u}) \, dx \leq \mathbf{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y}),$$

where

$$\begin{aligned} \mathbf{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y}) = & \\ & \frac{1 + \beta}{2} \int_{\Omega} (\nabla \mathbf{v} - \mathbf{A}^{-1} \mathbf{y}) \cdot (\mathbf{A} \nabla \mathbf{v} - \mathbf{y}) \, dx + \frac{1 + \beta}{2\beta} \frac{c_{\Omega}^2}{c_1^2} \|\operatorname{div} \mathbf{y} - \mathbf{f}\|^2 \quad (92) \end{aligned}$$

Certainly, the above estimate is applicable for the case $\mathbf{f} \in \mathbf{L}^2(\Omega)$ so that

$$\langle \mathbf{f}, \mathbf{w} \rangle = \int_{\Omega} \mathbf{f} \mathbf{w} \, dx,$$

and for $\mathbf{y} \in \mathbf{H}(\Omega, \mathbf{div})$.

Let $\{\mathbf{Y}_k^*\}$ be finite-dimensional subspaces of \mathbf{Y}^* such that

$$\begin{aligned} \mathbf{Y}_k^* &\in \mathbf{H}(\Omega, \mathbf{div}) \quad \text{for all } k = 1, 2, \dots; \\ \dim \mathbf{Y}_k^* &\rightarrow +\infty \quad \text{as } k \rightarrow \infty. \end{aligned}$$

We obtain computable upper bounds

$$\begin{aligned} M_{\oplus}^k = \inf_{\substack{\mathbf{y} \in \mathbf{Y}_k^* \\ \beta \in \mathbb{R}_+}} \left\{ \frac{1+\beta}{2} \int_{\Omega} (\nabla \mathbf{v} - \mathbf{A}^{-1} \mathbf{y}) \cdot (\mathbf{A} \nabla \mathbf{v} - \mathbf{y}) \, dx + \right. \\ \left. + \frac{1+\beta}{2\beta} \frac{c_{\Omega}^2}{c_1} \|\mathbf{div} \mathbf{y} - \mathbf{f}\|_{\Omega}^2 \right\}. \quad (93) \end{aligned}$$

Lower estimates

We have

$$\frac{1}{2} \int_{\Omega} \mathbf{A} \nabla(\mathbf{v} - \mathbf{u}) \cdot \nabla(\mathbf{v} - \mathbf{u}) \, dx \geq \mathbf{M}_{\Theta}(\mathbf{v}, \mathbf{w}), \quad \forall \mathbf{w} \in \mathbf{V}_0,$$

where

$$\mathbf{M}_{\Theta}(\mathbf{v}, \mathbf{w}) = -\frac{1}{2} \int_{\Omega} \mathbf{A} \nabla \mathbf{w} \cdot \nabla \mathbf{w} \, dx - \int_{\Omega} \mathbf{A} \nabla \mathbf{v} \cdot \nabla \mathbf{w} \, dx - \langle \mathbf{f}, \mathbf{w} \rangle.$$

Let $\{\mathbf{V}_{0\mathbf{k}}\}$ be finite-dimensional subspaces such that

$$\mathbf{V}_{0\mathbf{k}}^* \in \mathbf{V}_0 \quad \text{for all } \mathbf{k} = 1, 2, \dots;$$

$$\dim \mathbf{V}_{0\mathbf{k}} \rightarrow +\infty \quad \text{as } \mathbf{k} \rightarrow \infty.$$

Find the numbers

$$\mathbf{M}_{\Theta}^{\mathbf{k}} = \sup_{\mathbf{w}_{\mathbf{k}} \in \mathbf{V}_{0\mathbf{k}}} \mathbf{M}_{\Theta}(\mathbf{v}, \mathbf{w}_{\mathbf{k}}). \quad (94)$$

Both sequences \mathbf{M}_{\ominus}^k and \mathbf{M}_{\oplus}^k tend to $\frac{1}{2} \|\mathbf{v} - \mathbf{u}\|^2$ as $\mathbf{k} \rightarrow \infty$, provided that $\{\mathbf{Y}_{\mathbf{k}}^*\}$ and $\{\mathbf{V}_{0\mathbf{k}}\}$ possess necessary approximation properties (limit density). Note that if \mathbf{v} is a Galerkin approximation computed on $\mathbf{V}_{0\mathbf{k}}$, then $\mathbf{M}_{\ominus}(\mathbf{v}, \mathbf{w}_{\mathbf{k}}) = \mathbf{0}$. This means that to obtain a sensible lower estimate in this case, one must always use a finite-dimensional subspace that is *larger* than $\mathbf{V}_{0\mathbf{k}}$.

Neumann boundary condition

Consider the Neumann boundary condition

$$\boldsymbol{\nu} \cdot \mathbf{A} \nabla \mathbf{u} + \mathbf{F} = \mathbf{0} \quad \text{on} \quad \partial\Omega, \quad (95)$$

where $\boldsymbol{\nu}$ is the vector of unit outward normal to $\partial\Omega$. To apply the general scheme we set

$$\mathbf{V}_0 := \left\{ \mathbf{v} \in \mathbf{H}^1(\Omega) \mid \int_{\Omega} \mathbf{v} \, d\mathbf{x} = \mathbf{0} \right\}$$

and define $\boldsymbol{\Lambda}^* \mathbf{y} \in \mathbf{V}_0^*$ by the relation

$$\langle \boldsymbol{\Lambda}^* \mathbf{y}, \mathbf{w} \rangle = \int_{\Omega} \mathbf{y} \cdot \nabla \mathbf{w} \, d\mathbf{x}, \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

If \mathbf{y} is sufficiently regular then

$$\langle \boldsymbol{\Lambda}^* \mathbf{y}, \mathbf{w} \rangle = \int_{\Omega} (-\operatorname{div} \mathbf{y}) \mathbf{w} \, d\mathbf{x} + \int_{\partial\Omega} (\mathbf{y} \cdot \boldsymbol{\nu}) \mathbf{w} \, d\mathbf{x}.$$

Therefore, in such a case

$$\Lambda^* \mathbf{y} = [-\operatorname{div} \mathbf{y} \mid_{\Omega}; \quad \mathbf{y} \cdot \boldsymbol{\nu} \mid_{\partial\Omega}]$$

Also, we assume that \mathbf{F} and \mathbf{f} satisfy the equilibrium condition

$$\int_{\Omega} \mathbf{f} \, d\mathbf{x} + \int_{\partial\Omega} \mathbf{F} \, d\mathbf{x} = \mathbf{0}.$$

Assume that $\mathbf{f} \in \mathbf{L}^2(\Omega)$ and $\mathbf{F} \in \mathbf{L}^2(\partial\Omega)$. Then the Neumann problem has a solution defined by the integral identity

$$\int_{\Omega} \mathbf{A} \nabla \mathbf{u} \cdot \nabla \mathbf{w} \, d\mathbf{x} + \langle \ell, \mathbf{w} \rangle = \mathbf{0}, \quad \forall \mathbf{w} \in \mathbf{V}_0,$$

where

$$\langle \ell, \mathbf{w} \rangle = \int_{\Omega} \mathbf{f} \mathbf{w} \, d\mathbf{x} + \int_{\partial\Omega} \mathbf{F} \mathbf{w} \, d\mathbf{s}.$$

In general, $[\ell + \Lambda^* \mathbf{y}]$ is estimated in terms of the norms

$$\|\operatorname{div} \mathbf{y} - \mathbf{f}\|_{\mathbf{H}^{-1}} \quad \text{and} \quad \|\mathbf{y} \cdot \boldsymbol{\nu} + \mathbf{F}\|_{\mathbf{H}^{-1/2}}.$$

However, if we assume that \mathbf{y} possesses a certain regularity, so that

$$\mathbf{y} \in \mathbf{Q}^*(\Omega) := \{\mathbf{y} \in \mathbf{Y}^* \mid \operatorname{div} \mathbf{y} \in \mathbf{L}^2(\Omega), \mathbf{y} \cdot \boldsymbol{\nu} \in \mathbf{L}^2(\partial\Omega)\},$$

then

$$\langle \ell + \Lambda^* \mathbf{y}, \mathbf{w} \rangle = \int_{\Omega} (\mathbf{f} - \operatorname{div} \mathbf{y}) \mathbf{w} \, dx + \int_{\partial\Omega} (\mathbf{F} + \mathbf{y} \cdot \boldsymbol{\nu}) \mathbf{w} \, ds$$

and, therefore,

$$\begin{aligned} |\langle \ell + \Lambda^* \mathbf{y}, \mathbf{w} \rangle| &\leq \\ &\leq \|\operatorname{div} \mathbf{y} - \mathbf{f}\|_{2,\Omega} \|\mathbf{w}\|_{2,\Omega} + \|\mathbf{y} \cdot \boldsymbol{\nu} + \mathbf{F}\|_{2,\partial\Omega} \|\mathbf{w}\|_{2,\partial\Omega}. \end{aligned} \quad (96)$$

Let the constant \mathbf{c}_Ω be defined as

$$\frac{1}{\mathbf{c}_{(\Omega, \partial\Omega)}^2} = \inf_{\mathbf{w} \in \mathbf{V}_0} \frac{\int_{\Omega} \mathbf{A} \nabla \mathbf{w} \cdot \nabla \mathbf{w} \, dx}{\|\mathbf{w}\|_{2, \Omega}^2 + \|\mathbf{w}\|_{2, \partial\Omega}^2}.$$

Since the trace operator is bounded, this constant is finite. Therefore, (96) implies the estimate

$$\begin{aligned} |\langle \ell + \Lambda^* \mathbf{y}, \mathbf{w} \rangle| &\leq \\ &\leq \mathbf{c}_{(\Omega, \partial\Omega)} \left(\|\mathbf{div} \mathbf{y} - \mathbf{f}\|_{2, \Omega}^2 + \|\mathbf{y} \cdot \boldsymbol{\nu} + \mathbf{F}\|_{2, \partial\Omega}^2 \right)^{1/2} \|\Lambda \mathbf{w}\|^2 \end{aligned}$$

and the second term of the majorant is calculated as follows:

$$\begin{aligned} \mathbf{I} \ell + \Lambda^* \mathbf{y} \mathbf{I} &= \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\langle \ell + \Lambda^* \mathbf{y}, \mathbf{w} \rangle}{\|\Lambda \mathbf{w}\|} \leq \\ &\leq \mathbf{c}_{(\Omega, \partial\Omega)} \left(\|\mathbf{div} \mathbf{y} - \mathbf{f}\|_{2, \Omega}^2 + \|\mathbf{y} \cdot \boldsymbol{\nu} + \mathbf{F}\|_{2, \partial\Omega}^2 \right)^{1/2}. \end{aligned}$$

The term $\mathbf{D}(\Lambda \mathbf{v}, \mathbf{y})$ is defined as in the Dirichlét problem.

We see that the Majorants \mathbf{M}_{\oplus} for the two main boundary-value problems have different values of \mathbf{c}_{Ω} . In addition, the Neumann problem majorant contains an extra term

$$\|\mathbf{y} \cdot \boldsymbol{\nu} + \mathbf{F}\|_{2, \partial\Omega}$$

that penalizes violations of the Neumann boundary condition.

It is worth noting that if the given \mathbf{F} can be exactly reproduced by $\mathbf{y} \cdot \boldsymbol{\nu}$ for \mathbf{y} in a certain finite dimensional subspace $\mathbf{Y}_{\mathbf{k}}^*$, then one can compute $\mathbf{M}_{\oplus}^{\mathbf{k}}$ as

$$\mathbf{M}_{\oplus}^{\mathbf{k}} = \inf_{\substack{\mathbf{y} \in \mathbf{Y}_{\mathbf{k}}^*, \mathbf{y} \cdot \boldsymbol{\nu} = \mathbf{F} \text{ on } \partial\Omega \\ \beta \in \mathbb{R}_+}} \left\{ \frac{1+\beta}{2} \int_{\Omega} (\nabla \mathbf{v} - \mathbf{A}^{-1} \mathbf{y}) \cdot (\mathbf{A} \nabla \mathbf{v} - \mathbf{y}) \, \mathbf{d}\mathbf{x} + \right. \\ \left. + \frac{1+\beta}{2\beta} \mathbf{c}_{(\Omega, \partial\Omega)}^2 \|\operatorname{div} \mathbf{y} - \mathbf{f}\|_{\Omega}^2 \right\}. \quad (97)$$

Mixed boundary conditions

Let $\partial\Omega$ consist of two measurable nonintersecting parts $\partial_1\Omega$ and $\partial_2\Omega$, on which different boundary conditions are given:

$$\mathbf{u} = \mathbf{u}_0 \quad \text{on} \quad \partial_1\Omega,$$

$$\boldsymbol{\nu} \cdot \mathbf{A}\nabla\mathbf{u} + \mathbf{F} = \mathbf{0} \quad \text{on} \quad \partial_2\Omega.$$

Set

$$\mathbf{V}_0 := \left\{ \mathbf{v} \in \mathbf{H}^1(\Omega) \mid \mathbf{v} = \mathbf{0} \quad \text{on} \quad \partial_1\Omega \right\}$$

and

$$\langle \boldsymbol{\Lambda}^* \mathbf{y}, \mathbf{w} \rangle = \int_{\Omega} \mathbf{y} \cdot \nabla \mathbf{w} \, dx, \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

Assume that

$$\mathbf{f} \in \mathbf{L}^2(\Omega), \quad \mathbf{F} \in \mathbf{L}^2(\partial_2\Omega).$$

and \mathbf{y} possesses an extra regularity, namely,

$$\mathbf{y} \in \mathbf{Q}^*(\Omega) := \left\{ \mathbf{y} \in \mathbf{Y}^* \mid \operatorname{div} \mathbf{y} \in \mathbf{L}^2(\Omega), \mathbf{y} \cdot \boldsymbol{\nu} \in \mathbf{L}^2(\partial_2\Omega) \right\}.$$

Then, for any $\mathbf{w} \in \mathbf{V}_0$, we have

$$\langle \ell + \boldsymbol{\Lambda}^* \mathbf{y}, \mathbf{w} \rangle = \int_{\Omega} (\operatorname{div} \mathbf{y} - \mathbf{f}) \mathbf{w} \, dx + \int_{\partial_2\Omega} (\mathbf{y} \cdot \boldsymbol{\nu} + \mathbf{F}) \mathbf{w} \, ds,$$

Note that $\mathbf{p} \in \mathbf{Q}^*(\Omega)$!

Now, we obtain

$$|\langle \ell + \mathbf{\Lambda}^* \mathbf{y}, \mathbf{w} \rangle| \leq \| \mathbf{div} \mathbf{y} - \mathbf{f} \|_{2,\Omega} \| \mathbf{w} \|_{2,\Omega} + \| \mathbf{y} \cdot \boldsymbol{\nu} + \mathbf{F} \|_{2,\partial_2 \Omega} \| \mathbf{w} \|_{2,\partial_2 \Omega}.$$

Let γ and γ_* be two numbers such that $\gamma > \mathbf{1}$, $\gamma_* > \mathbf{1}$, $\frac{1}{\gamma} + \frac{1}{\gamma_*} = \mathbf{1}$.
Use the algebraic inequality

$$\mathbf{ab} + \mathbf{cd} \leq \sqrt{\gamma \mathbf{a}^2 + \gamma_* \mathbf{c}^2} \sqrt{\frac{1}{\gamma} \mathbf{b}^2 + \frac{1}{\gamma_*} \mathbf{d}^2}.$$

Then

$$|\langle \ell + \mathbf{\Lambda}^* \mathbf{y}, \mathbf{w} \rangle| \leq \left(\gamma \| \mathbf{div} \mathbf{y} - \mathbf{f} \|_{2,\Omega}^2 + \gamma_* \| \mathbf{y} \cdot \boldsymbol{\nu} + \mathbf{F} \|_{2,\partial_2 \Omega}^2 \right)^{1/2} \times \\ \times \left(\frac{1}{\gamma} \| \mathbf{w} \|_{2,\Omega}^2 + \frac{1}{\gamma_*} \| \mathbf{w} \|_{2,\partial_2 \Omega}^2 \right)^{1/2}.$$

Since (Friederichs type inequality)

$$\|\mathbf{w}\|_{2,\Omega}^2 \leq \mathbf{C}_F^2(\Omega) \|\nabla \mathbf{w}\|_{2,\Omega}^2, \quad \forall \mathbf{w} \in \mathbf{V}_0,$$

and (trace inequality)

$$\|\mathbf{w}\|_{2,\partial_2\Omega}^2 \leq \mathbf{C}_{tr}^2(\Omega, \partial_2\Omega) \|\mathbf{w}\|_{1,2,\Omega}^2, \quad \forall \mathbf{w} \in \mathbf{V}_0,$$

we find that

$$\begin{aligned} \frac{1}{\gamma} \|\mathbf{w}\|_{2,\Omega}^2 + \frac{1}{\gamma_*} \|\mathbf{w}\|_{2,\partial_2\Omega}^2 &\leq \\ &\leq \mathbf{C}_F^2 \frac{1}{\gamma} \|\nabla \mathbf{w}\|^2 + \mathbf{C}_{tr}^2 \frac{1}{\gamma_*} \left(\|\mathbf{w}\|_{2,\Omega}^2 + \|\nabla \mathbf{w}\|_{2,\Omega}^2 \right) \leq \\ &\leq \left(\mathbf{C}_F^2 \frac{1}{\gamma} + \mathbf{C}_{tr}^2 \frac{1}{\gamma_*} \left(1 + \mathbf{C}_F^2 \right) \right) \|\nabla \mathbf{w}\|_{2,\Omega}^2. \end{aligned}$$

Therefore, there exist a positive constant \mathbf{C}_γ such that

$$\frac{1}{\mathbf{C}_\gamma^2} = \inf_{\mathbf{w} \in \mathbf{V}_0} \frac{\int_{\Omega} \mathbf{A} \nabla \mathbf{w} \cdot \nabla \mathbf{w} \, dx}{\frac{1}{\gamma} \|\mathbf{w}\|_{2,\Omega}^2 + \frac{1}{\gamma_*} \|\mathbf{w}\|_{2,\partial_2\Omega}^2}.$$

The value of this constant can be estimated numerically by minimizing the above quotient on a sufficiently representative finite dimensional subspace. Besides, if \mathbf{C}_F and \mathbf{C}_{tr} are estimated, then

$$\mathbf{C}_\gamma^2 \leq \widehat{\mathbf{C}}_\gamma^2 := \left(\mathbf{C}_F^2 \frac{1}{\gamma} + \mathbf{C}_{tr}^2 (\mathbf{1} + \mathbf{C}_F^2) \frac{1}{\gamma_*} \right) \mathbf{c}_1^{-1},$$

so that an upper bound of \mathbf{C}_γ is directly computed. Now,

$$\begin{aligned} |\langle \ell + \mathbf{\Lambda}^* \mathbf{y}, \mathbf{w} \rangle| &\leq \\ &\leq \widehat{\mathbf{C}}_\gamma \left(\gamma \|\mathbf{div} \mathbf{y} - \mathbf{f}\|_{2,\Omega}^2 + \gamma_* \|\mathbf{y} \cdot \boldsymbol{\nu} + \mathbf{F}\|_{2,\partial_2\Omega}^2 \right)^{1/2} \|\nabla \mathbf{w}\|. \end{aligned}$$

From this estimate, we obtain

$$\|\ell + \Lambda^* \mathbf{y}\|^2 \leq \widehat{\mathbf{C}}_\gamma^2 \left(\gamma \|\operatorname{div} \mathbf{y} - \mathbf{f}\|_{2,\Omega}^2 + \frac{\gamma}{\gamma-1} \|\mathbf{y} \cdot \boldsymbol{\nu} + \mathbf{F}\|_{2,\partial_2\Omega}^2 \right).$$

Consider first the case, in which we simply set $\gamma = \gamma^* = 2$. Then

$$\widehat{\mathbf{C}}_{(\gamma=2)}^2 := \widehat{\mathbf{C}}_2^2 = \frac{1}{2} \left(\mathbf{C}_F^2 + \mathbf{C}_{\operatorname{tr}}^2 (\mathbf{1} + \mathbf{C}_F^2) \right) \mathbf{c}_1^{-1},$$

$$\|\ell + \Lambda^* \mathbf{y}\|^2 \leq 2\widehat{\mathbf{C}}_2^2 \left(\|\operatorname{div} \mathbf{y} - \mathbf{f}\|_{2,\Omega}^2 + \|\mathbf{y} \cdot \boldsymbol{\nu} + \mathbf{F}\|_{2,\partial_2\Omega}^2 \right).$$

and we find that

$$\begin{aligned} \mathbf{M}_\oplus(\mathbf{v}, \beta, \mathbf{y}) &= \frac{1+\beta}{2} \int_{\Omega} (\nabla \mathbf{v} - \mathbf{A}^{-1} \mathbf{y}) \cdot (\mathbf{A} \nabla \mathbf{v} - \mathbf{y}) \, dx + \\ &\quad + \frac{1+\beta}{2\beta} \widehat{\mathbf{C}}_2^2 \left(\|\operatorname{div} \mathbf{y} - \mathbf{f}\|_{2,\Omega}^2 + \|\mathbf{y} \cdot \boldsymbol{\nu} + \mathbf{F}\|_{2,\partial_2\Omega}^2 \right). \quad (98) \end{aligned}$$

This Majorant gives an upper bound of the deviation for any $\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0$, $\mathbf{y} \in \mathbf{Q}^*$, and $\beta > 0$.

M_{\oplus} for mixed boundary conditions

A more exact estimate is obtained if we define γ by minimizing of the Majorant. Then, we obtain

$$\begin{aligned} \mathbf{I} \ell + \mathbf{\Lambda}^* \mathbf{y} \mathbf{I}^2 &\leq \left(\mathbf{C}_F \|\mathbf{div} \mathbf{y} - \mathbf{f}\|_{2,\Omega} + \right. \\ &\quad \left. + \mathbf{C}_{\text{tr}} (\mathbf{1} + \mathbf{C}_F^2)^{1/2} \|\mathbf{y} \cdot \boldsymbol{\nu} + \mathbf{F}\|_{2,\partial_2\Omega} \right)^2 \mathbf{c}_1^{-2} \end{aligned}$$

and

$$\begin{aligned} \mathbf{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y}) &= \frac{1+\beta}{2} \int_{\Omega} (\nabla \mathbf{v} - \mathbf{A}^{-1} \mathbf{y}) \cdot (\mathbf{A} \nabla \mathbf{v} - \mathbf{y}) \, \mathbf{d}\mathbf{x} + \\ &\quad + \frac{1+\beta}{2\beta} \left(\mathbf{C}_F \|\mathbf{div} \mathbf{y} - \mathbf{f}\|_{2,\Omega} + \mathbf{C}_{\text{tr}} (\mathbf{1} + \mathbf{C}_F^2)^{1/2} \|\mathbf{y} \cdot \boldsymbol{\nu} + \mathbf{F}\|_{2,\partial_2\Omega} \right)^2 \mathbf{c}_1^{-2}. \quad (99) \end{aligned}$$

Majorant vanishes if and only if $\mathbf{v} = \mathbf{u}$ and $\mathbf{y} = \mathbf{A} \nabla \mathbf{u}$, it is continuous with respect to the convergence of \mathbf{v} in \mathbf{V} and \mathbf{y} in \mathbf{Q} .

Lower estimates

Lower estimates for the problems considered follow from the general ones obtained in the previous lecture. They have the form

$$\frac{1}{2} \int_{\Omega} \mathbf{A} \nabla(\mathbf{v} - \mathbf{u}) \cdot \nabla(\mathbf{v} - \mathbf{u}) \, d\mathbf{x} \geq \mathbf{M}_{\Theta}(\mathbf{v}, \mathbf{w}), \quad \forall \mathbf{w} \in \mathbf{V}_0,$$

where

$$\mathbf{M}_{\Theta}(\mathbf{v}, \mathbf{w}) = -\frac{1}{2} \int_{\Omega} \mathbf{A} \nabla(\mathbf{w} - \mathbf{v}) \cdot \nabla \mathbf{w} \, d\mathbf{x} - \int_{\Omega} \mathbf{f} \mathbf{w} \, d\mathbf{x} - \int_{\partial_2 \Omega} \mathbf{F} \mathbf{w} \, d\mathbf{s}.$$

Here \mathbf{V}_0 depends on the type of boundary conditions, and the integral over $\partial_2 \Omega$ must be eliminated in the case of Dirichlét problem.

Linear elasticity

Classical statement. The classical formulation is as follows:

Find a tensor-valued function $\boldsymbol{\sigma}^*$ (stress) and a vector-valued function \mathbf{u} (displacement) that satisfy the system of equations

$$\boldsymbol{\sigma}^* = \mathbb{L}\boldsymbol{\varepsilon}(\mathbf{u}) \quad \text{in} \quad \Omega, \quad (\text{Hooke's law})$$

$$\operatorname{div} \boldsymbol{\sigma}^* = \mathbf{f} \quad \text{in} \quad \Omega, \quad (\text{Equilibrium equation})$$

$$\mathbf{u} = \mathbf{u}_0 \quad \text{on} \quad \partial_1 \Omega,$$

$$\boldsymbol{\sigma}^* \boldsymbol{\nu} + \mathbf{F} = \mathbf{0} \quad \text{on} \quad \partial_2 \Omega.$$

where $\boldsymbol{\varepsilon}(\mathbf{u})$ is a symmetric part of the tensor $\nabla \mathbf{u}$.

Here Ω is a bounded domain with Lipschitz boundary $\partial\Omega$ that consists of two disjoint parts $\partial_1\Omega$ and $\partial_2\Omega$, $|\partial_1\Omega| > 0$, \mathbf{f} and \mathbf{F} are given forces and $\mathbb{L} = \{L_{ijkl}\}$ is the tensor of elasticity constants, which is subject to the conditions

$$\mathbf{C}_1|\boldsymbol{\varepsilon}|^2 \leq \mathbb{L}\boldsymbol{\varepsilon} : \boldsymbol{\varepsilon} \leq \mathbf{C}_2|\boldsymbol{\varepsilon}|^2, \quad \forall \boldsymbol{\varepsilon} \in \mathbb{M}_s^{n \times n},$$

and

$$L_{ijkl} = L_{jikm} = L_{kmij}, \quad L_{ijkl} \in \mathbf{L}^\infty(\Omega).$$

Generalized solution

Let

$$\mathbf{f} \in \mathbf{L}^2(\Omega, \mathbb{R}^n), \quad \mathbf{F} \in \mathbf{L}^2(\partial_2\Omega, \mathbb{R}^n).$$

Then, a generalized solution $\mathbf{u} \in \mathbf{V}_0 + \mathbf{u}_0$ is defined by the identity

$$\int_{\Omega} \mathbb{L} \varepsilon(\mathbf{u}) : \varepsilon(\mathbf{w}) \, d\mathbf{x} + \langle \ell, \mathbf{w} \rangle = 0, \quad \forall \mathbf{w} \in \mathbf{V}_0, \quad (100)$$

where

$$\langle \ell, \mathbf{w} \rangle = \int_{\Omega} \mathbf{f} \cdot \mathbf{w} \, d\mathbf{x} + \int_{\partial_2\Omega} \mathbf{F} \cdot \mathbf{w} \, ds.$$

Assume that \mathbf{u} is a smooth function and it satisfies the identity

$$\int_{\Omega} \mathbb{L} \varepsilon(\mathbf{u}) : \varepsilon(\mathbf{w}) \, d\mathbf{x} + \langle \ell, \mathbf{w} \rangle = \mathbf{0}, \quad \forall \mathbf{w} \in \mathbf{V}_0,$$

Then,

$$\int_{\Omega} (\mathbf{f} - \mathbf{div}(\mathbb{L} \varepsilon(\mathbf{u})) \cdot \mathbf{w} \, d\mathbf{x} + \int_{\partial_2 \Omega} ((\mathbb{L} \varepsilon(\mathbf{u})) \nu + \mathbf{F}) \cdot \mathbf{w} \, d\mathbf{s} = \mathbf{0},$$

$$\forall \mathbf{w} \in \mathbf{V}_0,$$

and we observe that in such a case the equilibrium equation and the Neumann boundary condition are satisfied in the classical sense.

Variational formulation

Note that the relation (100) is the Euler's equation for the functional

$$\mathbf{J}(\mathbf{v}) = \frac{1}{2} \int_{\Omega} \mathbb{L}\varepsilon(\mathbf{v}) : \varepsilon(\mathbf{v}) \, d\mathbf{x} + \langle \ell, \mathbf{v} \rangle .$$

Therefore, the respective boundary-value problem may be considered as a minimization problem for $\mathbf{J}(\mathbf{v})$ on the set

$$\mathbf{V}_0 := \{ \mathbf{v} \in \mathbf{H}^1(\Omega, \mathbb{R}^n) \mid \mathbf{v} = \mathbf{u}_0 \text{ on } \partial_1 \Omega \} .$$

To prove existence of a minimizer we must show the coercivity of $\mathbf{J}(\mathbf{v})$ on \mathbf{V}_0 . The key role in this belongs to the so-called Korn's inequality.

In the Dirichlet problem

$$\begin{aligned}
 \mathbf{J}(\mathbf{v}) &= \frac{1}{2} \int_{\Omega} \mathbb{L}\boldsymbol{\varepsilon}(\mathbf{v}) : \boldsymbol{\varepsilon}(\mathbf{v}) \, d\mathbf{x} + \langle \ell, \mathbf{v} \rangle \geq \\
 &\geq \frac{\mathbf{C}_1}{2} \|\boldsymbol{\varepsilon}(\mathbf{v})\|^2 - \|\mathbf{f}\| \|\mathbf{v}\| = \\
 &= \frac{\mathbf{C}_1}{2} \|\boldsymbol{\varepsilon}(\mathbf{u}_0 + \mathbf{w})\|^2 - \|\mathbf{f}\| \|\mathbf{u}_0 + \mathbf{w}\| \geq \\
 &\geq \frac{\mathbf{C}_1}{2} (\|\boldsymbol{\varepsilon}(\mathbf{u}_0)\| - \|\boldsymbol{\varepsilon}(\mathbf{w})\|)^2 - \|\mathbf{f}\| \|\mathbf{u}_0\| - \|\mathbf{f}\| \|\mathbf{w}\|.
 \end{aligned}$$

Thus, if we can prove that

$$\|\boldsymbol{\varepsilon}(\mathbf{w})\| \geq \mathbf{c} \|\nabla \mathbf{w}\| \quad \forall \mathbf{w} \in \mathring{\mathbf{H}}^1(\Omega),$$

then we would establish the coercivity of \mathbf{J} .

Korn's inequality

This inequality is required in various aspects of the mathematical analysis of elasticity problems. In the general form it states the equivalence of two norms:

$$\|\mathbf{w}\|_{1,2,\Omega} := \left(\int_{\Omega} (|\nabla \mathbf{w}|^2 + |\mathbf{w}|^2) \, d\mathbf{x} \right)^{1/2},$$

and

$$\|\mathbf{w}\|_{1,2,\Omega} := \left(\int_{\Omega} (|\varepsilon(\mathbf{w})|^2 + |\mathbf{w}|^2) \, d\mathbf{x} \right)^{1/2}.$$

Korn's inequality in $\mathring{\mathbf{H}}^1$

For the functions in $\mathring{\mathbf{H}}^1(\Omega)$ this fact is not difficult to prove. Indeed, for smooth functions

$$\begin{aligned} \int_{\Omega} |\varepsilon(\mathbf{w})|^2 \mathbf{d}\mathbf{x} &= \frac{1}{2} \|\nabla \mathbf{w}\|^2 + \frac{1}{2} \int_{\Omega} \sum_{ij} \mathbf{w}_{i,j} \mathbf{w}_{j,i} \mathbf{d}\mathbf{x} = \\ &= \frac{1}{2} \|\nabla \mathbf{w}\|^2 - \frac{1}{2} \int_{\Omega} \sum_{ij} \mathbf{w}_i \mathbf{w}_{j,ij} \mathbf{d}\mathbf{x} = \frac{1}{2} \|\nabla \mathbf{w}\|^2 + \frac{1}{2} \int_{\Omega} \sum_{ij} \mathbf{w}_{i,i} \mathbf{w}_{j,j} \mathbf{d}\mathbf{x} = \\ &= \frac{1}{2} \|\nabla \mathbf{w}\|^2 + \frac{1}{2} \int_{\Omega} \sum_i |\mathbf{w}_{i,i}|^2 \mathbf{d}\mathbf{x}. \end{aligned}$$

Since smooth functions are dense in $\mathring{\mathbf{H}}^1$, we find that

$$\|\nabla \mathbf{w}\| \leq \sqrt{2} \|\varepsilon(\mathbf{w})\| \quad \forall \mathbf{w} \in \mathring{\mathbf{H}}^1(\Omega). \quad (101)$$

By (101) we prove that the energy functional of the elasticity problem for the case of Dirichlet boundary conditions is coercive, i.e.,

$$\mathbf{J}(\mathbf{v}_k) \rightarrow +\infty, \quad \text{as } \|\nabla \mathbf{v}_k\| \rightarrow +\infty.$$

Rigid deflections

In the analysis of elasticity problems one more notion is often required. It is the so-called *Space of Rigid Deflections* that we denote $\mathbf{RD}(\Omega)$. This space is the kernel of the operator $\varepsilon(\mathbf{w})$, i.e. it contains vector-valued functions \mathbf{w} such that

$$\varepsilon(\mathbf{w}) = \mathbf{0}.$$

It can be defined as follows:

$$\mathbf{RD}(\Omega) := \{\mathbf{w} = \mathbf{w}_0 + \omega_0 \mathbf{x} \mid \mathbf{w}_0 \in \mathbb{R}^n, \omega_0 \in \mathbb{M}^{n \times n}\},$$

where $\omega_0(\mathbf{w}) = \frac{1}{2}(\nabla \mathbf{w} - (\nabla \mathbf{w})^T)$ is a skew-symmetric tensor associated with "rigid rotations".

Implications of the Korn's inequality

Theorem

Let Ω be a Lipschitz domain and $\partial_1\Omega$ is a nonempty connected part of the boundary. Then,

$$\|\mathbf{u}\|_{1,p,\Omega} \leq \mathbf{C} \left(\int_{\Omega} |\varepsilon(\mathbf{u})|^p \, d\mathbf{x} \right)^{\frac{1}{p}} \quad \forall \mathbf{u} \in \mathbf{V}_0, \quad \mathbf{p} \in (1, 2] \quad (102)$$

Proof. Assume the opposite. Then, for any $m \in \mathbb{N}$ we can find $\mathbf{v}^{(m)}$ such that $\mathbf{v}^{(m)} \in \mathbf{V}_0$ and

$$\|\mathbf{v}^{(m)}\|_{1,p,\Omega} > m \left(\int_{\Omega} |\varepsilon(\mathbf{v}^{(m)})|^p \, d\mathbf{x} \right)^{\frac{1}{p}}.$$

Set $\mathbf{w}^{(m)} = \frac{\mathbf{v}^{(m)}}{\|\mathbf{v}^{(m)}\|_{1,p,\Omega}}$, then

$$\|\mathbf{w}^{(m)}\|_{1,p,\Omega} = \mathbf{1} \quad \text{and} \quad \frac{1}{\mathbf{m}} \geq \left(\int_{\Omega} |\varepsilon(\mathbf{w}^{(m)})|^p \, d\mathbf{x} \right)^{\frac{1}{p}}.$$

Therefore,

$$\mathbf{w}^{(m)} \rightharpoonup \mathbf{w} \quad \text{in} \quad \mathbf{W}_p^1(\Omega, \mathbf{R}^n),$$

$$\mathbf{w}^{(m)} \rightarrow \mathbf{w} \quad \in \quad \mathbf{L}^p(\Omega, \mathbf{R}^n),$$

$$\|\varepsilon(\mathbf{w}^{(m)})\|_{p,\Omega} \rightarrow \mathbf{0} \quad \text{in} \quad \mathbf{L}^p(\Omega, \mathbf{R}^n).$$

From here we conclude that $\varepsilon(\mathbf{w}) = \mathbf{0}$.

Indeed, by the fact that a norm is weakly lower semicontinuous, we have

$$\mathbf{0} = \liminf_m \|\varepsilon(\mathbf{w}^{(m)})\|_{\mathbf{p},\Omega} \geq \|\varepsilon(\mathbf{w})\|_{\mathbf{p},\Omega}.$$

Thus, $\mathbf{w} \in \mathbf{RD}(\Omega) \cap \mathbf{V}$. There is only one such a function: $\mathbf{w} = \mathbf{0}$. It means that $\mathbf{w}^{(m)} \rightarrow \mathbf{0}$ in $\mathbf{L}^{\mathbf{p}}$. Now, we apply Korn's inequality

$$\|\mathbf{w}(\mathbf{m})\|_{1,\mathbf{p},\Omega} \leq \mathbf{C} \left(\int_{\Omega} (|\varepsilon(\mathbf{w}^{(m)})|^{\mathbf{p}} + |\mathbf{w}^{(m)}|^{\mathbf{p}}) \, \mathbf{d}\mathbf{x} \right)^{\frac{1}{\mathbf{p}}} \xrightarrow{\mathbf{m} \rightarrow \infty} \mathbf{0},$$

which shows that $\|\mathbf{w}(\mathbf{m})\|_{1,\mathbf{p},\Omega}$ tends to zero. But for any \mathbf{m} $\|\mathbf{w}^{(m)}\|_{\mathbf{p},1,\Omega} = \mathbf{1}$, so that such a behavior is impossible. We have arrived at a contradiction that proves the Theorem.

Another similar result is required for the Neumann problem. Define the set

$$\mathbf{V} = \left\{ \mathbf{v} \in \mathbf{W}_p^1(\Omega) \mid \int_{\Omega} \mathbf{v} \cdot \mathbf{w} \, d\mathbf{x} = 0 \quad \forall \mathbf{w} \in \mathbf{RD}(\Omega) \right\}.$$

Theorem

Let Ω be a bounded domain with Lipschitz boundary $\partial\Omega$. Then

$$\|\mathbf{u}\|_{1,p,\Omega} \leq \mathbf{C} \left(\int_{\Omega} |\varepsilon(\mathbf{u})|^p \, d\mathbf{x} \right)^{\frac{1}{p}} \quad \forall \mathbf{u} \in \mathbf{V}. \quad (103)$$

Proof. By the same arguments as before, we obtain a sequence $\mathbf{w}^{(m)} \in \mathbf{V}$ such that

$$\begin{aligned}\mathbf{w}^{(m)} &\rightharpoonup \mathbf{w} \quad \text{in} \quad \mathbf{W}_p^1(\Omega, \mathbf{R}^n), \\ \mathbf{w}^{(m)} &\rightarrow \mathbf{w} \quad \in \quad \mathbf{L}^p(\Omega, \mathbf{R}^n), \\ \|\varepsilon(\mathbf{w}^{(m)})\|_{p,\Omega} &\rightarrow \mathbf{0} \quad \text{in} \quad \mathbf{L}^p(\Omega, \mathbf{R}^n).\end{aligned}$$

By the arguments similar to those in the previous Theorem, we find that $\varepsilon(\mathbf{w}) = \mathbf{0}$ and, thus, $\mathbf{w} \in \mathbf{RD}(\Omega)$. In addition, for any $\bar{\mathbf{w}} \in \mathbf{RD}$, we have

$$\mathbf{0} = \int_{\Omega} \mathbf{w}^{(m)} \cdot \bar{\mathbf{w}} \, d\mathbf{x} = \int_{\Omega} \mathbf{w} \cdot \bar{\mathbf{w}} \, d\mathbf{x}.$$

But $\mathbf{w} \in \mathbf{RD}$, so that $\|\mathbf{w}\| = \mathbf{0}$, and by applying Korn's inequality we prove that $\|\mathbf{w}^{(m)}\|_{1,p,\Omega}$ tends to zero, what leads to a contradiction.

Estimates of deviations

Let \mathbf{v} and \mathbf{y} be some approximations of \mathbf{u} and $\boldsymbol{\sigma}^*$. Estimates of $\mathbf{v} - \mathbf{u}$ and $\mathbf{y} - \boldsymbol{\sigma}^*$ follow from the general scheme if we set

$$\begin{aligned} \mathbf{U} &= \mathbf{L}^2(\Omega, \mathbb{M}_s^{n \times n}), & \mathbf{V} &= \mathbf{H}^1(\Omega, \mathbb{R}^n), \\ \mathbf{V}_0 &= \{\mathbf{w} \in \mathbf{V} \mid \mathbf{w} = \mathbf{0} \text{ on } \partial_1 \Omega\}, \\ \|\mathbf{y}\|^2 &= \int_{\Omega} \mathbb{L} \mathbf{y} : \mathbf{y} \, dx, & \|\mathbf{y}\|_*^2 &= \int_{\Omega} \mathbb{L}^{-1} \mathbf{y} : \mathbf{y} \, dx, \end{aligned}$$

and $\boldsymbol{\Lambda} \mathbf{v} = \boldsymbol{\varepsilon}(\mathbf{v}) := \frac{1}{2} (\nabla \mathbf{v} + (\nabla \mathbf{v})^T)$. In this case,

$$\langle \boldsymbol{\Lambda}^* \mathbf{y}, \mathbf{w} \rangle = \int_{\Omega} \mathbf{y} : \boldsymbol{\varepsilon}(\mathbf{w}) \, dx, \quad \forall \mathbf{w} \in \mathbf{V}_0,$$

Now \mathbf{y} is a tensor-valued function and $\mathbf{y}\boldsymbol{\nu} = \mathbf{y}_{ij}\nu_j$ is a vector-function defined on $\partial\Omega$.

If

$$\mathbf{y} \in \mathbf{Q}^* := \{\mathbf{y} \in \mathbf{Y}^* \mid \mathbf{div} \mathbf{y} \in \mathbf{L}^2(\Omega, \mathbb{M}^{n \times n}), \mathbf{y}\boldsymbol{\nu} \in \mathbf{L}^2(\partial_2\Omega, \mathbb{R}^n)\}.$$

then

$$\langle \Lambda^* \mathbf{y}, \mathbf{w} \rangle = - \int_{\Omega} \mathbf{div} \mathbf{y} \cdot \mathbf{w} \, dx + \int_{\partial_2\Omega} (\mathbf{y}\boldsymbol{\nu}) \cdot \mathbf{w} \, d\Gamma$$

so that

$$\Lambda^* \mathbf{y} = \{-\mathbf{div} \mathbf{y} \mid_{\Omega}, \quad (\mathbf{y}\boldsymbol{\nu}) \mid_{\partial_2\Omega}\}.$$

Upper estimates

By applying the general estimate, we obtain the following upper estimate:

$$\frac{1}{2} \int_{\Omega} \mathbb{L} \varepsilon(\mathbf{v} - \mathbf{u}) : \varepsilon(\mathbf{v} - \mathbf{u}) \, \mathbf{d}\mathbf{x} \leq \mathbf{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y}),$$

where

$$\mathbf{M}_{\oplus}(\mathbf{v}, \beta, \mathbf{y}) = \frac{1+\beta}{2} \mathbf{D}(\varepsilon\mathbf{v}, \mathbf{y}) + \frac{1+\beta}{2\beta} \mathbf{I} \, \boldsymbol{\Lambda}^* \mathbf{y} + \ell \mathbf{I}^2$$

and

$$\begin{aligned} \mathbf{D}(\varepsilon(\mathbf{v}), \mathbf{y}) &= \frac{1}{2} \int_{\Omega} \left(\mathbb{L} \varepsilon(\mathbf{v}) : \varepsilon(\mathbf{v}) \mathbb{L}^{-1} \mathbf{y} : \mathbf{y} - 2\varepsilon(\mathbf{v}) : \mathbf{y} \right) \, \mathbf{d}\mathbf{x} = \\ &= \int_{\Omega} (\varepsilon(\mathbf{u}) - \mathbb{L}^{-1} \mathbf{y}) : (\mathbb{L} \varepsilon(\mathbf{u}) - \mathbf{y}) \, \mathbf{d}\mathbf{x}. \end{aligned}$$

If $\mathbf{y} \in \mathbf{Q}^*$, then

$$\begin{aligned}
 \mathbf{I} \Lambda^* \mathbf{y} + \ell \mathbf{I} &= \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\langle \Lambda^* \mathbf{y} + \ell, \mathbf{w} \rangle}{\|\Lambda \mathbf{w}\|} = \\
 &= \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\int_{\Omega} (\mathbf{y} : \boldsymbol{\varepsilon}(\mathbf{w}) + \mathbf{f} \cdot \mathbf{w}) \, dx + \int_{\partial_2 \Omega} \mathbf{F} \cdot \mathbf{w} \, ds}{\|\boldsymbol{\varepsilon}(\mathbf{w})\|} = \\
 &= \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\int_{\Omega} (\mathbf{f} - \operatorname{div} \mathbf{y}) \cdot \mathbf{w} \, dx + \int_{\partial_2 \Omega} (\mathbf{F} + \mathbf{y} \nu) \cdot \mathbf{w} \, ds}{\|\boldsymbol{\varepsilon}(\mathbf{w})\|} \leq \\
 &\leq \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\|\mathbf{f} - \operatorname{div} \mathbf{y}\|_{2, \Omega} \|\mathbf{w}\|_{2, \Omega} + \|\mathbf{F} + \mathbf{y} \nu\|_{\partial_2 \Omega} \|\mathbf{w}\|_{\partial_2 \Omega}}{\|\boldsymbol{\varepsilon}(\mathbf{w})\|}.
 \end{aligned}$$

Let \mathbf{C}_Ω be a constant in the inequality

$$\int_{\Omega} |\mathbf{w}|^2 \, d\mathbf{x} + \int_{\partial_2 \Omega} |\mathbf{w}|^2 \, d\mathbf{s} \leq \mathbf{C}_\Omega^2 \|\boldsymbol{\varepsilon}(\mathbf{w})\|_\Omega^2, \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

Note that the existence of such a constant follows from the Korn's inequality. Indeed, the inequality

$$\int_{\Omega} |\mathbf{w}|^2 \, d\mathbf{x} + \int_{\partial_2 \Omega} |\mathbf{w}|^2 \, d\mathbf{s} \leq \hat{\mathbf{C}}_\Omega^2 \|\nabla(\mathbf{w})\|_\Omega^2, \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

for the tensor–gradient $\nabla(\mathbf{w})$ follows from the Friederichs type inequality for the vector–valued functions and the respective trace theorems. By (102) we recall that for the functions in \mathbf{V}_0

$$\|\nabla(\mathbf{w})\|_\Omega \leq \mathbf{C} \|\boldsymbol{\varepsilon}(\mathbf{w})\|_\Omega$$

with a certain constant \mathbf{C} and the estimate follows.

In practice, values of \mathbf{C}_Ω can be estimated by minimizing the quotient

$$\frac{\|\varepsilon(\mathbf{w})\|_\Omega^2}{\int_\Omega |\mathbf{w}|^2 \, d\mathbf{x} + \int_{\partial_2 \Omega} |\mathbf{w}|^2 \, d\mathbf{s}}$$

over sufficiently representative finite dimensional space $\mathbf{V}_{0h} \subset \mathbf{V}_0$.

Let us now return to finding an upper bound of the quantity $\|\Lambda^* \mathbf{y} + \ell\|$.

By the inequality $ab + cd \leq \sqrt{a^2 + c^2} \sqrt{b^2 + d^2}$, we obtain

$$\begin{aligned} \|\Lambda^* \mathbf{y} + \ell\| &\leq \\ &\leq \left(\|\operatorname{div} \mathbf{y} - \mathbf{f}\|_\Omega^2 + \|\mathbf{F} + \mathbf{y}\nu\|_{\partial_2 \Omega}^2 \right)^{1/2} \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{(\|\mathbf{w}\|_\Omega^2 + \|\mathbf{w}\|_{\partial_2 \Omega}^2)^{1/2}}{\|\varepsilon(\mathbf{w})\|} \leq \\ &\leq \mathbf{C}_\Omega \mathbf{c}_1^{-1/2} \left(\|\operatorname{div} \mathbf{y} - \mathbf{f}\|_\Omega^2 + \|\mathbf{F} + \mathbf{y}\nu\|_{\partial_2 \Omega}^2 \right)^{1/2} \end{aligned}$$

Error Majorant for mixed boundary conditions

Hence, we arrive at the Majorant \mathbf{M}_{\oplus} :

$$\begin{aligned} \mathbf{M}_{\oplus}(\boldsymbol{\varepsilon}(\mathbf{v}), \mathbf{y}) = & \frac{1+\beta}{2} \int_{\Omega} (\boldsymbol{\varepsilon}(\mathbf{v}) - \mathbb{L}^{-1}\mathbf{y}) : (\mathbb{L}\boldsymbol{\varepsilon}(\mathbf{v}) - \mathbf{y}) \, d\mathbf{x} + \\ & + \frac{1+\beta}{2\beta c_1} \mathbf{C}_{\Omega}^2 \left(\|\operatorname{div} \mathbf{y} - \mathbf{f}\|_{\Omega}^2 + \|\mathbf{F} + \mathbf{y}\nu\|_{\partial_2\Omega}^2 \right). \quad (104) \end{aligned}$$

It has a clear physical meaning. The first term of \mathbf{M}_{\oplus} is nonnegative and vanishes if and only if

$$\mathbf{y} = \mathbb{L}\boldsymbol{\varepsilon}(\mathbf{v}).$$

It penalizes violations of the **Hooke's law**. The meaning of the second term is obvious: it contains \mathbf{L}^2 -norms of other two relations, which gives errors in the **equilibrium equation** and **boundary condition** for the stress tensor.

Thus, the majorant not only gives an idea of the overall value of the error, but also shows its physically sensible parts.

Let $\{\mathbf{Y}_k^*\} \subset \mathbf{H}^1(\Omega, \mathbb{M}^{n \times n})$ be a collection of finite-dimensional subspaces that satisfy the limit density condition. Then, (104) generates a sequence of computable upper bounds

$$M_{\oplus}^k = \inf_{\substack{\mathbf{y} \in \mathbf{Y}_k^* \\ \beta \in \mathbb{R}_+}} \left\{ \frac{1+\beta}{2} \int_{\Omega} (\mathbb{L} \varepsilon(\mathbf{v}) : \varepsilon(\mathbf{v}) + \mathbb{L}^{-1} \mathbf{y} : \mathbf{y} - 2\varepsilon(\mathbf{v}) : \mathbf{y}) dx + \frac{1+\beta}{2\beta c_1} \mathbf{C}_{\Omega}^2 (\|\mathbf{div} \mathbf{y} - \mathbf{f}\|_{\Omega}^2 + \|\mathbf{F} + \mathbf{y}\nu\|_{\partial_2 \Omega}^2) \right\},$$

which tends to the exact value of the error.

Lower estimates

Lower estimates also follow from the general theory. We have

$$\frac{1}{2} \int_{\Omega} \mathbb{L} \varepsilon(\mathbf{v} - \mathbf{u}) : \varepsilon(\mathbf{v} - \mathbf{u}) \, dx \geq \mathbf{M}_{\Theta}(\mathbf{v}, \mathbf{w}), \quad \forall \mathbf{w} \in \mathbf{V}_0,$$

where

$$\begin{aligned} \mathbf{M}_{\Theta}(\mathbf{v}, \mathbf{w}) = & -\frac{1}{2} \int_{\Omega} \mathbb{L} \varepsilon(\mathbf{w}) : \varepsilon(\mathbf{w}) \, dx - \int_{\Omega} \mathbb{L} \varepsilon(\mathbf{v}) : \varepsilon(\mathbf{w}) \, dx - \\ & - \int_{\Omega} \mathbf{f} \cdot \mathbf{w} \, dx - \int_{\partial_2 \Omega} \mathbf{F} \cdot \mathbf{w} \, ds. \end{aligned}$$

By the same arguments as for the diffusion equation one can prove that

$$\frac{1}{2} \int_{\Omega} \mathbb{L} \varepsilon(\mathbf{v} - \mathbf{u}) : \varepsilon(\mathbf{v} - \mathbf{u}) \, dx = \sup_{\mathbf{w} \in \mathbf{V}_0} \mathbf{M}_{\Theta}(\mathbf{v}, \mathbf{w}_k).$$

By the maximization the functional \mathbf{M}_Θ on a sequence of finite-dimensional spaces $\mathbf{V}_{0k} \subset \mathbf{V}_0$, we obtain a sequence of computable lower bounds

$$\mathbf{M}_\Theta^k = \sup_{\mathbf{w} \in \mathbf{V}_{0k}} \mathbf{M}_\Theta(\mathbf{v}, \mathbf{w}_k).$$

If the spaces \mathbf{V}_{0k} satisfy the limit density condition stated, then the sequence of numbers $\{\mathbf{M}_\Theta\}$ tends to $\frac{1}{2} \|\varepsilon(\mathbf{v} - \mathbf{u})\|^2$.

Linear elliptic equations of the fourth order

Now, we consider the problem

$$\nabla \cdot \nabla \cdot (\mathbf{B} \nabla \nabla \mathbf{u}) = \mathbf{f} \quad \text{in } \Omega, \quad (105)$$

$$\mathbf{u} = \frac{\partial \mathbf{u}}{\partial \nu} = \mathbf{0} \quad \text{on } \partial \Omega. \quad (106)$$

Here $\Omega \subset \mathbb{R}^2$, ν denotes the outward unit normal to the boundary, and $\mathbf{B} = \{\mathbf{b}_{ijkl}\} \in \mathcal{L}(\mathbb{M}_s^{2 \times 2}, \mathbb{M}_s^{2 \times 2})$. We assume that $\mathbf{b}_{ijkl} = \mathbf{b}_{jikl} = \mathbf{b}_{klij}$,

$$\alpha_1 |\boldsymbol{\eta}|^2 \leq \mathbf{B} \boldsymbol{\eta} : \boldsymbol{\eta} \leq \alpha_2 |\boldsymbol{\eta}|^2, \quad \forall \boldsymbol{\eta} \in \mathbb{M}_s^{2 \times 2},$$

and

$$\mathbf{f} \in \mathbf{L}^2(\Omega), \quad \mathbf{b}_{ijkl} \in \mathbf{L}^\infty(\Omega).$$

To apply the general scheme, we set

$$\begin{aligned}\mathbf{U} &= \mathbf{L}^2(\Omega, \mathbb{M}_s^{2 \times 2}), & \mathbf{V} &= \mathbf{H}^2(\Omega), \\ \mathbf{V}_0 &= \{\mathbf{w} \in \mathbf{V} \mid \mathbf{w} = \frac{\partial \mathbf{w}}{\partial \nu} = \mathbf{0} \text{ on } \partial\Omega\},\end{aligned}$$

and define \mathbf{A} as the Hessian operator. Now, the basic integral identity has the form

$$\int_{\Omega} \mathbf{B} \nabla \nabla \mathbf{u} : \nabla \nabla \mathbf{w} \, d\mathbf{x} = \int_{\Omega} \mathbf{f} \mathbf{w} \, d\mathbf{x} \quad \forall \mathbf{w} \in \mathbf{V}_0. \quad (107)$$

By \mathbf{B}^{-1} we denote the inverse tensor, which satisfies the double inequality

$$\alpha_2^{-1} |\eta|^2 \leq \mathbf{B}^{-1} \eta : \eta \leq \alpha_1^{-1} |\eta|^2, \quad \forall \eta \in \mathbb{M}_s^{2 \times 2},$$

The spaces \mathbf{Y} and \mathbf{Y}^* are equipped with norms

$$\| \mathbf{y} \|^2 = \int_{\Omega} \mathbf{B} \mathbf{y} : \mathbf{y} \, dx; \quad \| \mathbf{y} \|_*^2 = \int_{\Omega} \mathbf{B}^{-1} \mathbf{y} : \mathbf{y} \, dx,$$

$$\langle \ell, \mathbf{w} \rangle = - \int_{\Omega} \mathbf{f} \mathbf{w} \, dx,$$

and

$$\mathbf{Q}_\ell^* = \{ \mathbf{y} \in \mathbf{Y}^* \mid \int_{\Omega} \mathbf{y} : \nabla \nabla \mathbf{w} \, dx = \int_{\Omega} \mathbf{f} \mathbf{w} \, dx, \quad \forall \mathbf{w} \in \mathbf{V}_0 \}.$$

Since

$$\| \nabla \nabla \mathbf{w} \| \geq \alpha_3 \| \mathbf{w} \|_{2,2,\Omega} \quad \forall \mathbf{w} \in \mathbf{V}_0,$$

we have the required version of the coercivity condition

$$\| \Lambda \mathbf{w} \| \geq \mathbf{c}_3 \| \mathbf{w} \|_{\mathbf{V}}.$$

Problem (105) and (106) is associated with two variational problems.

Problem \mathcal{P} . Find $\mathbf{u} \in \mathbf{V}_0$ such that

$$\mathbf{J}(\mathbf{u}) = \inf_{\mathbf{v} \in \mathbf{V}_0} \mathbf{J}(\mathbf{v}),$$

where

$$\mathbf{J}(\mathbf{v}) = \frac{1}{2} \int_{\Omega} \mathbf{B} \nabla \nabla \mathbf{v} : \nabla \nabla \mathbf{v} \, dx - \int_{\Omega} \mathbf{f} \mathbf{w} \, dx.$$

Problem \mathcal{P}^* . Find $\mathbf{p} \in \mathbf{Q}_{\ell}^*$ such that

$$\mathbf{I}^*(\mathbf{p}) = \sup_{\mathbf{q} \in \mathbf{Q}_{\ell}^*} \mathbf{I}^*(\mathbf{q}),$$

where

$$\mathbf{I}^*(\mathbf{q}) = -\frac{1}{2} \int_{\Omega} \mathbf{B}^{-1} \mathbf{q} : \mathbf{q} \, dx.$$

Two basic relations for the deviations that we derived in the previous Lecture now come in the form:

$$\| \nabla \nabla (\mathbf{v} - \mathbf{u}) \|^2 + \| \mathbf{q} - \mathbf{p} \|^2_* = 2(\mathbf{J}(\mathbf{v}) - \mathbf{I}^*(\mathbf{q})), \quad (108)$$

and

$$\begin{aligned} \| \nabla \nabla (\mathbf{v} - \mathbf{u}) \|^2 + \| \mathbf{q} - \mathbf{p} \|^2_* &= 2\mathbf{D}(\nabla \nabla \mathbf{v}, \mathbf{q}) = \\ &= \int_{\Omega} \left(\mathbf{B} \nabla \nabla \mathbf{v} : \nabla \nabla \mathbf{v} + \mathbf{B}^{-1} \mathbf{q} : \mathbf{q} - 2 \nabla \nabla \mathbf{v} : \mathbf{q} \right) \mathbf{d}\mathbf{x}. \end{aligned} \quad (109)$$

They hold for any $\mathbf{v} \in \mathbf{V}_0$ and $\mathbf{q} \in \mathbf{Q}_\ell^*$.

Also, from the general theory it readily follows the first a posteriori estimate:

$$\frac{1}{2} \| \nabla \nabla (\mathbf{v} - \mathbf{u}) \|^2 \leq (\mathbf{1} + \beta) \mathbf{D}(\nabla \nabla \mathbf{v}, \mathbf{y}) + \left(\mathbf{1} + \frac{1}{\beta} \right) \frac{\mathbf{d}_\ell^2(\mathbf{y})}{2}, \quad (110)$$

where $\mathbf{d}_\ell^2(\mathbf{y}) = \inf_{\mathbf{q} \in \mathbf{Q}_\ell^*} \| \mathbf{q} - \mathbf{y} \|^2_*$.

Note that

$$\int_{\Omega} \mathbf{y} : \nabla \nabla \mathbf{w} \, d\mathbf{x} = \int_{\Omega} (\mathbf{div} \mathbf{div} \, \mathbf{y}) \mathbf{w} \, d\mathbf{x}, \quad \forall \mathbf{w} \in \mathbf{V}_0,$$

so that $\mathbf{\Lambda}^* : \mathbf{Y}^* \rightarrow \mathbf{H}^{-2}(\Omega)$ is the operator $\mathbf{div} \mathbf{div}$.

Next,

$$\langle \ell + \mathbf{\Lambda}^* \mathbf{y}, \mathbf{w} \rangle = \int_{\Omega} (\mathbf{y} : \nabla \nabla \mathbf{w} - \mathbf{f} \mathbf{w}) \, d\mathbf{x}$$

and, therefore,

$$\mathbf{d}_{\ell}^2(\mathbf{y}) = \mathbf{I} \ell + \mathbf{\Lambda}^* \mathbf{y} \mathbf{I} = \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\int_{\Omega} (\mathbf{y} : \nabla \nabla \mathbf{w} - \mathbf{f} \mathbf{w}) \, d\mathbf{x}}{\|\nabla \nabla \mathbf{w}\|}.$$

If

$$\mathbf{y} \in \mathbf{H}(\mathbf{divdiv}, \Omega) := \left\{ \mathbf{y} \in \mathbf{L}^2(\Omega, \mathbb{M}_s^{n \times n}) \mid \mathbf{divdiv} \mathbf{y} \in \mathbf{L}^2(\Omega) \right\},$$

then this quantity is estimated by the relation

$$\begin{aligned} \mathbf{I} \ell + \mathbf{\Lambda}^* \mathbf{y} \mathbf{I} &\leq \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\|\mathbf{divdiv} \mathbf{y} - \mathbf{f}\|_{\Omega} \|\mathbf{w}\|_{\Omega}}{\|\nabla \nabla \mathbf{w}\|} \leq \\ &\leq \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\|\mathbf{divdiv} \mathbf{y} - \mathbf{f}\|_{\Omega} \|\mathbf{w}\|_{\Omega}}{\alpha_1 \|\nabla \nabla \mathbf{w}\|} \leq \frac{\mathbf{C}_{1\Omega}}{\alpha_1} \|\mathbf{divdiv} \mathbf{y} - \mathbf{f}\|_{\Omega}, \end{aligned}$$

in which $\mathbf{C}_{1\Omega}$ is a constant in the inequality

$$\|\mathbf{w}\|_{\Omega} \leq \mathbf{C}_{1\Omega} \|\nabla \nabla \mathbf{w}\|_{\Omega} \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

Now, we obtain the first variant of a posteriori estimate for the biharmonic type problem.

First a posteriori estimate

$$\frac{1}{2} \|\nabla\nabla(\mathbf{v} - \mathbf{u})\|^2 \leq (1 + \beta) \mathbf{D}(\nabla\nabla\mathbf{v}, \mathbf{y}) + \left(1 + \frac{1}{\beta}\right) \frac{\mathbf{C}_{1\Omega}^2}{2\alpha_1^2} \|\mathbf{divdiv}\mathbf{y} - \mathbf{f}\|_{\Omega}^2, \quad (111)$$

Here, \mathbf{y} is an arbitrary tensor-valued function from $\mathbf{H}(\mathbf{div}\mathbf{div}, \Omega)$ and β is a positive real number. However, this is rather demanding in relation to the dual variable \mathbf{y} (which must have square summable \mathbf{divdiv}). To avoid technical difficulties that rises from this condition, we estimate the negative norm in a different way.

$$\begin{aligned}
\mathbf{I} \ell + \mathbf{\Lambda}^* \mathbf{y} \mathbf{I} &= \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\int_{\Omega} (\mathbf{y} : \nabla \nabla \mathbf{w} - \mathbf{f} \mathbf{w}) \, \mathrm{d}\mathbf{x}}{\|\nabla \nabla \mathbf{w}\|} = \\
&= \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\int_{\Omega} (\mathbf{y} : \nabla \nabla \mathbf{w} + \boldsymbol{\eta} \cdot \nabla \mathbf{w} + \operatorname{div} \boldsymbol{\eta} \mathbf{w} - \mathbf{f} \mathbf{w}) \, \mathrm{d}\mathbf{x}}{\|\nabla \nabla \mathbf{w}\|} = \\
&= \frac{\int_{\Omega} (-\operatorname{div} \boldsymbol{\eta} \cdot \nabla \mathbf{w} + \boldsymbol{\eta} \cdot \nabla \mathbf{w} + \operatorname{div} \boldsymbol{\eta} \mathbf{w} - \mathbf{f} \mathbf{w}) \, \mathrm{d}\mathbf{x}}{\|\nabla \nabla \mathbf{w}\|} \leq \\
&\leq \frac{\mathbf{C}_{2\Omega}}{\alpha_1} \|\operatorname{div} \mathbf{y} - \boldsymbol{\eta}\|_{\Omega} + \frac{\mathbf{C}_{1\Omega}}{\alpha_1} \|\operatorname{div} \boldsymbol{\eta} - \mathbf{f}\|_{\Omega}.
\end{aligned}$$

Here, $\boldsymbol{\eta}$ is an arbitrary vector-valued function from $\mathbf{H}(\operatorname{div}, \Omega)$ and $\mathbf{C}_{2\Omega}$ is a constant in the inequality

$$\|\nabla \mathbf{w}\|_{\Omega} \leq \mathbf{C}_{2\Omega} \|\nabla \nabla \mathbf{w}\|_{\Omega} \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

Second a posteriori estimate

Then, we arrive at the estimate

$$\frac{1}{2} \|\nabla\nabla(\mathbf{v} - \mathbf{u})\|^2 \leq (1 + \beta) \mathbf{D}(\nabla\nabla\mathbf{v}, \mathbf{y}) + \left(1 + \frac{1}{\beta}\right) \frac{1}{2\alpha_1^2} (\mathbf{C}_{2\Omega} \|\operatorname{div} \mathbf{y} - \boldsymbol{\eta}\|_{\Omega} + \mathbf{C}_{1\Omega} \|\operatorname{div} \boldsymbol{\eta} - \mathbf{f}\|_{\Omega})^2, \quad (112)$$

in which $\mathbf{y} \in \boldsymbol{\Sigma}_{\operatorname{div}}(\Omega)$ and $\boldsymbol{\eta} \in \mathbf{H}(\operatorname{div}, \Omega)$.

This estimate was obtained in

P. Neittaanmäki and S. Repin. A posteriori error estimates for boundary-value problems related to the biharmonic operator, *East-West J.Numer. Math.*, 9(2001)

Note that

$$\|\mathbf{w}\| \leq \mathbf{C}_F \|\nabla \mathbf{w}\|_{\Omega} \leq \mathbf{C}_F \mathbf{C}_{2\Omega} \|\nabla \nabla \mathbf{w}\|_{\Omega} \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

where \mathbf{C}_F is a constant in the Friederichs inequality. Therefore, $\mathbf{C}_{1\Omega} \leq \mathbf{C}_F \mathbf{C}_{2\Omega}$. In view of this, we obtain a slightly different form of the deviation estimate:

$$\begin{aligned} \frac{1}{2} \|\nabla \nabla (\mathbf{v} - \mathbf{u})\|^2 &\leq (1 + \beta) \mathbf{D}(\nabla \nabla \mathbf{v}, \mathbf{y}) + \\ &+ \left(1 + \frac{1}{\beta}\right) \frac{\mathbf{C}_{2\Omega}^2}{2\alpha_1^2} (\|\mathbf{div} \mathbf{y} - \boldsymbol{\eta}\|_{\Omega} + \mathbf{C}_F \|\mathbf{div} \boldsymbol{\eta} - \mathbf{f}\|_{\Omega})^2, \quad (113) \end{aligned}$$

For boundary conditions of other types, the deviation majorants can be derived by arguments similar to those used in Lecture 6.

Lower estimates of the deviation from \mathbf{u}

Lower estimates follow from the general estimate discussed in Lecture 5.

We have

$$\frac{1}{2} \|\nabla\nabla(\mathbf{v} - \mathbf{w})\|^2 \geq \mathbf{M}_\Theta(\mathbf{v}, \mathbf{w}) \quad \mathbf{w} \in \mathbf{V}_0, \quad (114)$$

where

$$\mathbf{M}_\Theta(\mathbf{v}, \mathbf{w}) := -\frac{1}{2} \|\nabla\nabla\mathbf{w}\|^2 - \int_{\Omega} (\mathbf{B}\nabla\nabla\mathbf{v} : \nabla\nabla\mathbf{w} - \mathbf{f}\mathbf{w}) dx.$$

Lecture 5

A POSTERIORI ESTIMATES IN NON-ENERGY QUANTITIES

If we can estimate the energy error norm, then we can estimate all other quantities subject to it.

General Principle

If

$$\| \mathbf{u} - \mathbf{v} \| \leq \mathbf{M}_\oplus$$

and

$$\phi(\mathbf{u} - \mathbf{v}) \leq \mathbf{C} \| \mathbf{u} - \mathbf{u}_a \|$$

then

$$\phi(\mathbf{u} - \mathbf{v}) \leq \mathbf{C} \mathbf{M}_\oplus$$

Example

By the embedding theorems we know that for any Lipschitz domain

If $q \geq 1$, and $1 + \frac{n}{q} \geq \frac{n}{2}$, then $\mathbf{W}^{1,2}(\Omega)$ is continuously embedded in $\mathbf{L}^q(\Omega)$, i.e.,

$$\|\mathbf{w}\|_{q,\Omega} \leq \mathbf{C}(\Omega, q) \|\mathbf{w}\|_{1,2,\Omega}$$

If $1 > \frac{n}{2}$, then $\mathbf{W}^{1,2}(\Omega)$ is compactly embedded in $\mathbf{C}(\overline{\Omega})$, i.e.,

$$\|\mathbf{w}\|_{\mathbf{C}(\Omega)} \leq \hat{\mathbf{C}}(\Omega) \|\mathbf{w}\|_{1,2,\Omega}$$

If the conditions of above theorems hold, then we obtain

$$\|\mathbf{u} - \mathbf{v}\|_{\mathbf{q}, \Omega} \leq \mathbf{C}(\Omega, \mathbf{q}) \mathbf{M}_{\oplus}(\mathbf{v}, \mathbf{y})$$

Note that if $\mathbf{n} = \mathbf{2}$ then the above relation holds for any finite \mathbf{q} .

$$\text{esssup}_{\mathbf{x} \in \Omega} \|\mathbf{u}(\mathbf{x}) - \mathbf{v}(\mathbf{x})\| \leq \mathbf{C}(\Omega, \mathbf{q}) \mathbf{M}_{\oplus}(\mathbf{v}, \mathbf{y})$$

ERROR ESTIMATES IN TERMS OF GOAL-ORIENTED FUNCTIONALS

$$\ell(\mathbf{u} - \mathbf{v}) \leq \mathbf{M}_{\ell \oplus}$$

Example:

$$\langle \ell, \mathbf{u} - \mathbf{v} \rangle = \int_{\Omega} \varphi(\mathbf{u} - \mathbf{v}) \, \mathbf{d}\mathbf{x},$$

where the weight φ is a locally supported function.

Let us consider the methods of goal-oriented error control on the paradigm of our Basic Linear Problem.

$$(\mathcal{A}\Lambda\mathbf{u}, \Lambda\mathbf{w}) = \langle \mathbf{f}, \mathbf{w} \rangle, \quad \forall \mathbf{w} \in \mathbf{V}_0. \quad (115)$$

$$\mathcal{A} \in \mathcal{L}(\mathbf{U}, \mathbf{U}),$$

$$\Lambda \in \mathcal{L}(\mathbf{V}, \mathbf{U}),$$

\mathbf{V} is a Banach space, \mathbf{V}_0 is a subspace of \mathbf{V} , $f \in \mathbf{V}_0^*$,

\mathbf{U} is a Hilbert space equipped with scalar product (\cdot, \cdot) ,

$\langle \cdot, \cdot \rangle$ denotes the duality pairing of \mathbf{V}_0 and \mathbf{V}_0^* .

We assume that the operators \mathcal{A} and Λ satisfy the conditions

$$\mathbf{c}_1 \|\mathbf{y}\|_{\mathbf{U}}^2 \leq (\mathcal{A}\mathbf{y}, \mathbf{y}) \leq \mathbf{c}_2 \|\mathbf{y}\|_{\mathbf{U}}^2, \forall \mathbf{y} \in \mathbf{U}, \quad (116)$$

$$\mathbf{c}_3 \|\mathbf{w}\|_{\mathbf{V}} \leq \|\Lambda \mathbf{w}\|_{\mathbf{U}}, \quad \forall \mathbf{w} \in \mathbf{V}_0, \quad (117)$$

Our aim is to derive estimates of the quantity

$$\langle \ell, \mathbf{u} - \mathbf{v} \rangle,$$

where \mathbf{v} is an arbitrary element from \mathbf{V}_0 .

A commonly accepted way of the derivation of such estimates is exposed in the works of

M. Ainsworth, R. Becker, C. Johnson,

K. Eriksson, P. Hansbo, T. Oden,

S. Prudome, R. Rannacher, E. Suli and other scientist, see, e.g.,

M. Ainsworth and J. T. Oden, *A posteriori error estimation in the finite element method*, Numer. Math., 60(1992), pp. 429-463.

W. Bangerth and R. Rannacher, *Adaptive finite element methods for differential equations*, Birkhäuser, Berlin, 2003.

It consists of the attraction the **adjoint problem**.

Adjoint problem

Find $\mathbf{u}_a \in \mathbf{V}_0$ such that

$$(\mathcal{A}^* \Lambda \mathbf{u}_a, \Lambda \mathbf{w}) = \langle \ell, \mathbf{w} \rangle \quad \forall \mathbf{w} \in \mathbf{V}_0, \quad (118)$$

where \mathcal{A}^* is the operator adjoint to \mathcal{A} , i.e.,

$$(\mathcal{A} \mathbf{y}, \mathbf{z}) = (\mathbf{y}, \mathcal{A}^* \mathbf{z}) \quad \forall \mathbf{y}, \mathbf{z} \in \mathbf{U}.$$

Proposition

Let \mathbf{u} , \mathbf{v} and \mathbf{u}_a , \mathbf{v}_a be exact and approximate solutions of the primal and adjoint problems respectively. Then,

$$\langle \ell, \mathbf{u} - \mathbf{v} \rangle := \mathbf{E}(\mathbf{v}, \mathbf{v}_a) = \mathbf{E}_0(\mathbf{v}, \mathbf{v}_a) + \mathbf{E}_1(\mathbf{v}, \mathbf{v}_a), \quad (119)$$

where

$$\mathbf{E}_0(\mathbf{v}, \mathbf{v}_a) = \langle \mathbf{f}, \mathbf{v}_a \rangle - (\mathcal{A}\mathbf{\Lambda}\mathbf{v}, \mathbf{\Lambda}\mathbf{v}_a) \quad (120)$$

and

$$\mathbf{E}_1(\mathbf{v}, \mathbf{v}_a) = (\mathcal{A}\mathbf{\Lambda}(\mathbf{u} - \mathbf{v}), \mathbf{\Lambda}(\mathbf{u}_a - \mathbf{v}_a)). \quad (121)$$

Proof.

Since \mathbf{u}_a is a solution of the adjoint problem, we have

$$\begin{aligned} \langle \ell, \mathbf{u} - \mathbf{v} \rangle &= (\mathcal{A}^* \boldsymbol{\Lambda} \mathbf{u}_a, \boldsymbol{\Lambda}(\mathbf{u} - \mathbf{v})) = \\ &= (\mathcal{A} \boldsymbol{\Lambda}(\mathbf{u} - \mathbf{v}), \boldsymbol{\Lambda} \mathbf{u}_a) = \\ &= (\mathcal{A} \boldsymbol{\Lambda}(\mathbf{u} - \mathbf{v}), \boldsymbol{\Lambda}(\mathbf{u}_a - \mathbf{v}_a)) + \langle \mathbf{f}, \mathbf{v}_a \rangle - (\mathcal{A} \boldsymbol{\Lambda} \mathbf{v}, \boldsymbol{\Lambda} \mathbf{v}_a). \end{aligned}$$



$\mathbf{E}_0(\mathbf{v}, \mathbf{v}_a)$ is explicitly computable !

$\mathbf{E}_1(\mathbf{v}, \mathbf{v}_a)$ contains unknown solutions of the two problems.

Let \mathbf{V}_h and \mathbf{V}_τ be two finite-dimensional subspaces of \mathbf{V}_0 , and let $\mathbf{v} = \mathbf{u}_h$ and $\mathbf{v}_a = \mathbf{u}_{a\tau}$, where \mathbf{u}_h and $\mathbf{u}_{a\tau}$ are solutions of the problems

$$(\mathcal{A}\Lambda\mathbf{u}_h, \Lambda\mathbf{w}_h) = \langle \mathbf{f}, \mathbf{w}_h \rangle \quad \forall \mathbf{w}_h \in \mathbf{V}_h, \quad (122)$$

$$(\mathcal{A}^*\Lambda\mathbf{u}_{a\tau}, \Lambda\mathbf{w}_\tau) = \langle \ell, \mathbf{w}_\tau \rangle \quad \forall \mathbf{w}_\tau \in \mathbf{V}_\tau. \quad (123)$$

In the particular case $\mathbf{V}_h \equiv \mathbf{V}_\tau$, the relation (122) implies that $\mathbf{E}_0(\mathbf{u}_h, \mathbf{u}_{a\tau}) = \mathbf{0}$ so that there remains the term containing the product of the (unknown) energy errors. On the contrary, for noncoinciding spaces both terms \mathbf{E}_0 and \mathbf{E}_1 are present. Moreover, it is easy to show that the term \mathbf{E}_0 dominates if \mathbf{v}_a is close to \mathbf{u}_a . Indeed, if $\mathbf{v}_a \rightarrow \mathbf{u}_a$ in \mathbf{V} , then $\Lambda(\mathbf{u}_a - \mathbf{v}_a) \rightarrow \mathbf{0}$ in \mathbf{U} , so that

$$\mathbf{E}_1(\mathbf{v}, \mathbf{v}_a) \rightarrow \mathbf{0}.$$

However,

$$\mathbf{E}_0(\mathbf{v}, \mathbf{v}_a) + \mathbf{E}_1(\mathbf{v}, \mathbf{v}_a) = \langle \ell, \mathbf{u} - \mathbf{v} \rangle \neq \mathbf{0}.$$

Hence, if \mathbf{v}_a is sufficiently close to \mathbf{u}_a , then the directly computable term $\mathbf{E}_0(\mathbf{v}, \mathbf{v}_a)$ contains the major part of the estimated quantity.

Key question: how to estimate

$$\mathbf{E}_1(\mathbf{v}, \mathbf{v}_a) = (\mathcal{A}\boldsymbol{\Lambda}(\mathbf{u} - \mathbf{v}), \boldsymbol{\Lambda}(\mathbf{u}_a - \mathbf{v}_a)).$$

Two ways.

- Economical method that gives correct presentation on the error but does not give guaranteed error bounds,
- Computationally expensive method that provides two-sided guaranteed error bounds for the goal-oriented functional.

Economical way: use post-processing methods

Main idea: recover the unknown functions \mathbf{Au} and \mathbf{Au}_a by some post-processing techniques.

See: S. Korotov, P. Neittaanmaki and S. Repin,
*A posteriori error estimation of goal-oriented quantities by the
superconvergence patch recovery*, J. Numer. Math. 11 (2003), 1, pp. 33-59.

Let G_h and G_τ be averaging operators defined on \mathbf{V}_h and \mathbf{V}_τ , respectively. Replace $\mathbf{E}(\mathbf{u}_h, \mathbf{u}_{a\tau})$ by the directly computable functional

$$\tilde{\mathbf{E}}(\mathbf{u}_h, \mathbf{u}_{a\tau}) := \mathbf{E}_0(\mathbf{u}_h, \mathbf{u}_{a\tau}) + \tilde{\mathbf{E}}_1(\mathbf{u}_h, \mathbf{u}_{a\tau}), \quad (124)$$

where

$$\tilde{\mathbf{E}}_1(\mathbf{u}_h, \mathbf{u}_{a\tau}) = (\mathcal{A}(G_h(\boldsymbol{\Lambda}\mathbf{u}_h) - \boldsymbol{\Lambda}\mathbf{u}_h), G_\tau(\boldsymbol{\Lambda}\mathbf{u}_{a\tau}) - \boldsymbol{\Lambda}\mathbf{u}_{a\tau}).$$

If the operators G_h and G_τ perform a proper recovery of $\boldsymbol{\Lambda}\mathbf{u}_h$ and $\boldsymbol{\Lambda}\mathbf{u}_{a\tau}$, then it is natural to expect that the difference between $\mathbf{E}_1(\mathbf{u}_h, \mathbf{u}_{a\tau})$ and $\tilde{\mathbf{E}}(\mathbf{u}_h, \mathbf{u}_{a\tau})$ is given by higher order terms and, thus, the latter quantity can successfully be used instead of \mathbf{E}_1 .

Two-sided estimates

Main idea: is to apply two-sided error estimates given by \mathbf{M}_\oplus and \mathbf{M}_\ominus to derive sharp bounds of $(\mathcal{A}\Lambda(\mathbf{u} - \mathbf{v}), \Lambda(\mathbf{u}_a - \mathbf{v}_a))$ using a special representation of this term.

See

P. Neittaanmäki and S. Repin, *Reliable methods for computer simulation, Error control and a posteriori estimates*, Elsevier, New York, 2004.

Certainly, we may simply use the estimate

$$(\mathcal{A}\boldsymbol{\Lambda}(\mathbf{u} - \mathbf{v}), \boldsymbol{\Lambda}(\mathbf{u}_a - \mathbf{v}_a)) \leq \| \boldsymbol{\Lambda}(\mathbf{u} - \mathbf{v}) \| \| \boldsymbol{\Lambda}(\mathbf{u}_a - \mathbf{v}_a) \| .$$

However, it may considerably **overestimate** the error.
Therefore, it is suggested to use the algebraic relation

$$2\mathbf{a}\mathbf{b} = \left(\alpha\mathbf{a} + \frac{1}{\alpha}\mathbf{b} \right)^2 - \alpha^2\mathbf{a}^2 - \frac{1}{\alpha^2}\mathbf{b}^2.$$

We have

$$\langle \ell, \mathbf{u} - \mathbf{v} \rangle := \mathbf{E}(\mathbf{v}, \mathbf{v}_a) = \mathbf{E}_0(\mathbf{v}, \mathbf{v}_a) + \mathbf{E}_1(\mathbf{v}, \mathbf{v}_a), \quad (125)$$

where \mathbf{v} is an approximation of \mathbf{u} , \mathbf{v}_a is an approximation of \mathbf{u}_a , and

$$\begin{aligned} \mathbf{E}_0(\mathbf{v}, \mathbf{v}_a) &= \langle \mathbf{f}, \mathbf{v}_a \rangle - (\mathcal{A}\mathbf{\Lambda}\mathbf{v}, \mathbf{\Lambda}\mathbf{v}_a), \\ \mathbf{E}_1(\mathbf{v}, \mathbf{v}_a) &= (\mathcal{A}\mathbf{\Lambda}\delta_f, \mathbf{\Lambda}\delta_\ell), \quad \delta_f = (\mathbf{u} - \mathbf{v}), \quad \delta_\ell = (\mathbf{u}_a - \mathbf{v}_a). \end{aligned}$$

By the algebraic identity

$$\begin{aligned} 2(\mathcal{A}\mathbf{\Lambda}\delta_\ell, \mathbf{\Lambda}\delta_f) &= \\ &= \|\mathbf{\Lambda}(\alpha\delta_f + \frac{1}{\alpha}\delta_\ell)\|^2 - \alpha^2 \|\mathbf{\Lambda}\delta_f\|^2 - \frac{1}{\alpha^2} \|\mathbf{\Lambda}\delta_\ell\|^2. \end{aligned} \quad (126)$$

Assume that for the primal and adjoint problems we have found good upper and lower error bounds in the energy norm, i.e., we have \mathbf{y}_f^* , \mathbf{y}_ℓ^* , β_f , β_ℓ , \mathbf{w}_f , and \mathbf{w}_ℓ such that

$$\mathbf{m}_f := \mathcal{M}_\ominus(\mathbf{v}, \mathbf{w}_f) \leq \frac{1}{2} \|\Lambda \delta_f\|^2 \leq \mathcal{M}_\oplus(\mathbf{v}, \beta_f, \mathbf{y}_f^*) := \mathbf{M}_f,$$
$$\mathbf{m}_\ell := \mathcal{M}_\ominus(\mathbf{v}_a, \mathbf{w}_\ell) \leq \frac{1}{2} \|\Lambda \delta_\ell\|^2 \leq \mathcal{M}_\oplus(\mathbf{v}_a, \beta_\ell, \mathbf{y}_\ell^*) := \mathbf{M}_\ell.$$

Note that the quantity

$$\| \mathbf{A}(\alpha\delta_{\mathbf{f}} + \frac{1}{\alpha}\delta_{\ell}) \|^2 = \| \mathbf{A}(\alpha\mathbf{u} + \frac{1}{\alpha}\mathbf{u}_{\mathbf{a}} - \alpha\mathbf{v} - \frac{1}{\alpha}\mathbf{v}_{\mathbf{a}}) \|^2$$

can be viewed as the norm of the difference between the exact solution $\mathbf{u}_{\mathbf{f}\ell}^{\alpha} \in \mathbf{V}_0 + (\alpha + \frac{1}{\alpha})\mathbf{u}_0$ of the problem

$$(\mathbf{A}\mathbf{L}\mathbf{u}_{\mathbf{f}\ell}^{\alpha}, \mathbf{L}\mathbf{w}) + \langle \alpha\mathbf{f} + \frac{1}{\alpha}\ell, \mathbf{w} \rangle = 0, \quad \forall \mathbf{w} \in \mathbf{V}_0$$

and the function $\mathbf{v}_{\mathbf{a}\ell}^{\alpha} = \alpha\mathbf{v} + \frac{1}{\alpha}\mathbf{v}_{\mathbf{a}} \in \mathbf{V}_0 + (\alpha + \frac{1}{\alpha})\mathbf{u}_0$, which is an approximation of $\mathbf{u}_{\mathbf{f}\ell}^{\alpha}$.

To obtain two-sided bounds of $\| \mathbf{\Lambda}(\alpha\delta_{\mathbf{f}} + \frac{1}{\alpha}\delta_{\ell}) \| ^2$ we use already known functions $\mathbf{y}_{\mathbf{f}}^*$, \mathbf{y}_{ℓ}^* , $\mathbf{w}_{\mathbf{f}}$, and \mathbf{w}_{ℓ} and compute the numbers

$$\mathbf{m}_{\mathbf{f}\ell} = \mathcal{M}_{\ominus}(\mathbf{u}_{\mathbf{a}_{\mathbf{f}\ell}}^{\alpha}, \alpha\mathbf{w}_{\mathbf{f}} + \frac{1}{\alpha}\mathbf{w}_{\ell}),$$

$$\mathbf{M}_{\mathbf{f}\ell} = \mathcal{M}_{\oplus}(\mathbf{u}_{\mathbf{a}_{\mathbf{f}\ell}}^{\alpha}, \beta, \alpha\mathbf{y}_{\mathbf{f}}^* + \frac{1}{\alpha}\mathbf{y}_{\ell}^*).$$

We note that $\mathbf{m}_{\mathbf{f}\ell} = \mathbf{m}_{\mathbf{f}\ell}(\alpha)$, $\mathbf{M}_{\mathbf{f}\ell} = \mathbf{M}_{\mathbf{f}\ell}(\alpha, \beta)$ and positive numbers α and β can be taken arbitrary. By (126), we obtain

$$(\mathcal{A}\mathbf{\Lambda}\delta_{\mathbf{f}}, \mathbf{\Lambda}\delta_{\ell}) \leq \inf_{\alpha, \beta \in \mathbb{R}_+} \{ \mathbf{M}_{\mathbf{f}\ell} - \alpha^2 \mathbf{m}_{\mathbf{f}} - \frac{\mathbf{m}_{\ell}}{\alpha^2} \} := \mathfrak{M}_{\oplus},$$

$$(\mathcal{A}\mathbf{\Lambda}\delta_{\mathbf{f}}, \mathbf{\Lambda}\delta_{\ell}) \geq \sup_{\alpha \in \mathbb{R}_+} \{ \mathbf{m}_{\mathbf{f}\ell} - \alpha^2 \mathbf{M}_{\mathbf{f}} - \frac{\mathbf{M}_{\ell}}{\alpha^2} \} := \mathfrak{M}_{\ominus}.$$

Recalling (125), we now deduce a two-sided estimate

$$\mathfrak{M}_{\ominus} + \mathbf{E}_0(\mathbf{v}, \mathbf{v}_a) \leq \langle \ell, \mathbf{u} - \mathbf{v} \rangle \leq \mathfrak{M}_{\oplus} + \mathbf{E}_0(\mathbf{v}, \mathbf{v}_a), \quad (127)$$

If \mathbf{y}_f^* , \mathbf{y}_ℓ^* , \mathbf{w}_f , and \mathbf{w}_ℓ provide accurate two-sided estimates of the energy error norms in the primal and adjoint problems, then $\mathbf{M}_{f\ell}$ and $\mathbf{m}_{f\ell}$ furnish accurate two-sided estimates for the norm $\| \mathbf{\Lambda}(\alpha\delta_f + \frac{1}{\alpha}\delta_\ell) \|^2$ and, consequently, (127) yields sharp upper and lower bounds of $\langle \ell, \mathbf{u} - \mathbf{v} \rangle$.

Error estimates in terms of seminorms

Estimates of $\mathbf{u} - \mathbf{v}$ computed on a set of linear functionals can be used for an evaluation of this difference in some seminorms.

Let ω be a subdomain in Ω and $\{\varphi_1, \varphi_2, \dots, \varphi_d\}$ be a set of functions that vanish in $\Omega \setminus \omega$.

By the method described above, we can estimate the quantities

$$\mathbf{l}_{\varphi_s} := \int_{\omega} \varphi_s(\mathbf{x})(\mathbf{u}(\mathbf{x}) - \mathbf{v}(\mathbf{x})) \, d\mathbf{x}.$$

Using these quantities we can estimate the error in terms of **local seminorms**

Let

$$\|\mathbf{u} - \mathbf{v}\|_{2,\omega}^2 := \int_{\omega} |\mathbf{u}(\mathbf{x}) - \mathbf{v}(\mathbf{x})|^2 \, d\mathbf{x}.$$

Note that

$$\|\mathbf{u} - \mathbf{v}\|_{2,\omega} = \sup_{\eta \in \mathbf{L}^2(\omega)} \frac{\int_{\omega} \eta(\mathbf{u} - \mathbf{v}) \, d\mathbf{x}}{\|\eta\|_{2,\omega}},$$

and the supremum is attained if $\eta = \mathbf{u} - \mathbf{v}$.

Thus, if the difference $\mathbf{u} - \mathbf{v}$ is known to belong to a certain set $\Upsilon(\omega) \subset \mathbf{L}^2(\omega)$, then the problem is reduced to the evaluation of the *seminorm*

$$|\mathbf{u} - \mathbf{v}|_{\Upsilon} := \sup_{\eta \in \Upsilon(\omega)} \frac{\int_{\omega} \eta(\mathbf{u} - \mathbf{v}) \mathbf{d}\mathbf{x}}{\|\eta\|_{2,\omega}}.$$

In general, a seminorm does not provide complete information on the error, because the relation $|\mathbf{u} - \mathbf{v}|_{\Upsilon} = \mathbf{0}$ does not mean that $\mathbf{u} = \mathbf{v}$. Nevertheless, seminorms may give a useful information if Υ contains sufficiently large amount of linearly independent functions.

Let Υ be a subspace made by φ_s , i.e.,

$$\Upsilon = \text{Span} \{ \varphi_1, \varphi_2, \dots, \varphi_d \} = \left\{ \sum_s \alpha_s \varphi_s \right\}.$$

Then,

$$|\mathbf{u} - \mathbf{v}|_{\Upsilon} = \sup_{\mathbf{a}_i \in \mathbb{R}} \frac{\sum_{i=1}^d \mathbf{a}_i \int_{\omega} \varphi_i (\mathbf{u} - \mathbf{v}) \mathbf{d}\mathbf{x}}{\left(\sum_{i,j=1}^d \mathbf{a}_i \mathbf{a}_j \int_{\omega} \varphi_i \varphi_j \mathbf{d}\mathbf{x} \right)^{\frac{1}{2}}}. \quad (128)$$

This problem is equivalent to the following one:

$$\inf_{\mathbf{a}_i \in \mathbb{R}} \sum_{i,j=1}^d \mathbf{a}_i \mathbf{a}_j \int_{\omega} \varphi_i \varphi_j \mathbf{d}\mathbf{x}.$$

for all \mathbf{a}_i such that

$$\sum_{i=1}^d \mathbf{a}_i \int_{\omega} \varphi_i (\mathbf{u} - \mathbf{v}) \mathbf{d}\mathbf{x} = \mathbf{1}$$

The latter problem can be reformulated as:

Problem. Find the quantity

$$\kappa^2(\omega, \Upsilon) = \inf_{\mathbf{a} \in \mathbf{Q}} \mathbf{B} \mathbf{a} \cdot \mathbf{a}, \quad (129)$$

where $\mathbf{a} = (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_d)$,

$$\mathbf{Q} = \{\mathbf{a} \in \mathbb{R}^d \mid \mathbf{a} \cdot \mathbf{l} = 1\}, \quad \mathbf{l} = (l_{\varphi_1}, l_{\varphi_2}, \dots, l_{\varphi_d}),$$

and

$$\mathbf{B} = \{\mathbf{b}_{ij}\}_{i,j=1,d}, \quad \mathbf{b}_{ij} = \int_{\omega} \varphi_i \varphi_j \, \mathbf{d}\mathbf{x}.$$

The problem (129) has a minimax form

$$\begin{aligned}\kappa^2(\omega, \boldsymbol{\Upsilon}) &= \inf_{\mathbf{a} \in \mathbb{R}^d} \sup_{\lambda \in \mathbb{R}} \{ \mathbf{B} \mathbf{a} \cdot \mathbf{a} + \lambda [(\mathbf{a} \cdot \mathbf{I}) - \mathbf{1}] \} \geq \\ &\geq \sup_{\lambda \in \mathbb{R}} \inf_{\mathbf{a} \in \mathbb{R}^d} \{ \mathbf{B} \mathbf{a} \cdot \mathbf{a} + \lambda [(\mathbf{a} \cdot \mathbf{I}) - \mathbf{1}] \}. \quad (130)\end{aligned}$$

Assume that the functions φ_i are chosen in such a way that \mathbf{B} is a positive definite matrix. Then (130) holds as equality. The minimization problem that stands on the right-hand side in (130) has a solution

$$\mathbf{a} = -\frac{\lambda}{2} \mathbf{B}^{-1} \mathbf{I}.$$

Therefore,

$$\kappa^2(\omega, \Upsilon) = \sup_{\lambda \in \mathbb{R}} \left\{ -\frac{\lambda^2}{4} \mathbf{B}^{-1} \mathbf{l} \cdot \mathbf{l} - \lambda \right\} = \frac{1}{\mathbf{B}^{-1} \mathbf{l} \cdot \mathbf{l}},$$

and

$$|\mathbf{u} - \mathbf{v}|_{\Upsilon} = \frac{1}{\kappa(\omega, \Upsilon)} = (\mathbf{B}^{-1} \mathbf{l} \cdot \mathbf{l})^{\frac{1}{2}}.$$

We see that the value of $|\mathbf{u} - \mathbf{v}|_{\Upsilon}$ is not difficult to find, provided that the functions $\{\varphi_s\}$ are properly chosen and the respective quantities l_{φ_s} are defined.

Let us show that if $\mathbf{v} \in \Upsilon(\omega)$ and \mathbf{u} is a sufficiently smooth function, then **seminorm $|\mathbf{u} - \mathbf{v}|_{\Upsilon}$ gives a good estimate of $\|\mathbf{u} - \mathbf{v}\|_{2,\omega}$** . Let $\Upsilon(\omega)$ be a subspace of $L^2(\omega)$ and $\hat{\mathbf{u}}$ be an arbitrary function from $\Upsilon(\omega)$. Then

$$\begin{aligned} \|\mathbf{u} - \mathbf{v}\|_{2,\omega} &= \sup_{\eta \in L^2(\omega)} \frac{\int_{\omega} ((\mathbf{u} - \hat{\mathbf{u}})\eta + (\hat{\mathbf{u}} - \mathbf{v})\eta) \, d\mathbf{x}}{\|\eta\|_{2,\omega}} \leq \\ &\leq \sup_{\eta \in L^2(\omega)} \frac{\int_{\omega} (\hat{\mathbf{u}} - \mathbf{v})\eta \, d\mathbf{x}}{\|\eta\|_{2,\omega}} + \|\mathbf{u} - \hat{\mathbf{u}}\|_{2,\omega}. \end{aligned}$$

Since $\hat{\mathbf{u}} - \mathbf{v} \in \Upsilon(\omega)$, we have

$$\|\mathbf{u} - \mathbf{v}\|_{2,\omega} \leq \sup_{\eta \in \Upsilon(\omega)} \frac{\int (\hat{\mathbf{u}} - \mathbf{v})\eta \, \mathbf{d}\mathbf{x}}{\|\eta\|_{2,\omega}} + \|\mathbf{u} - \hat{\mathbf{u}}\|_{2,\omega}.$$

By rearranging again the numerator of the fraction on the right-hand side, we arrive at the estimate

$$\begin{aligned} \|\mathbf{u} - \mathbf{v}\|_{2,\omega} &= \sup_{\eta \in \Upsilon(\omega)} \frac{\int ((\hat{\mathbf{u}} - \mathbf{u})\eta + (\mathbf{u} - \mathbf{v})\eta) \, \mathbf{d}\mathbf{x}}{\|\eta\|_{2,\omega}} + \|\mathbf{u} - \hat{\mathbf{u}}\|_{2,\omega} \leq \\ &\leq \sup_{\eta \in \Upsilon(\omega)} \frac{\int (\mathbf{u} - \mathbf{v})\eta \, \mathbf{d}\mathbf{x}}{\|\eta\|_{2,\omega}} + 2\|\mathbf{u} - \hat{\mathbf{u}}\|_{2,\omega} = |\mathbf{u} - \mathbf{v}|_{\Upsilon} + 2\|\mathbf{u} - \hat{\mathbf{u}}\|_{2,\omega}, \end{aligned}$$

in which $\hat{\mathbf{u}}$ is an arbitrary function from $\Upsilon(\omega)$.

Hence, we see that

$$\|\mathbf{u} - \mathbf{v}\|_{2,\omega} \leq |\mathbf{u} - \mathbf{v}|_{\Upsilon} + 2\text{dist}(\mathbf{u}, \Upsilon(\omega)), \quad (131)$$

where

$$\text{dist}(\mathbf{u}, \Upsilon(\omega)) = \inf_{\hat{\mathbf{u}} \in \Upsilon(\omega)} \|\mathbf{u} - \hat{\mathbf{u}}\|_{2,\omega}$$

is the distance between \mathbf{u} and $\Upsilon(\omega)$.

If $\Upsilon(\Omega)$ is a set of polynomials and \mathbf{u} is a smooth function then $\delta(\mathbf{u}, \Upsilon(\omega))$ can be estimated with the help of well-known results of approximation theory.

By similar arguments, we obtain

$$\begin{aligned} \|\mathbf{u} - \mathbf{v}\|_{2,\omega} &= \sup_{\eta \in L^2(\omega)} \frac{\int_{\omega} ((\mathbf{u} - \hat{\mathbf{u}})\eta + (\hat{\mathbf{u}} - \mathbf{v})\eta) \, d\mathbf{x}}{\|\eta\|_{2,\omega}} \geq \\ &\geq \sup_{\eta \in L^2(\omega)} \frac{\int_{\omega} (\hat{\mathbf{u}} - \mathbf{v})\eta \, d\mathbf{x}}{\|\eta\|_{2,\omega}} - \|\mathbf{u} - \hat{\mathbf{u}}\|_{2,\omega}. \end{aligned}$$

Since $\hat{\mathbf{u}} - \mathbf{v}$ belongs to $\Upsilon(\omega)$, we can replace in the first integral $\eta \in L^2$ by $\eta \in \Upsilon$.

Then, we obtain

$$\begin{aligned}
\|\mathbf{u} - \mathbf{v}\|_{2,\omega} &\geq \sup_{\eta \in \Upsilon(\omega)} \frac{\int ((\hat{\mathbf{u}} - \mathbf{u})\eta + (\mathbf{u} - \mathbf{v})\eta) \, d\mathbf{x}}{\|\eta\|_{2,\omega}} - \|\mathbf{u} - \hat{\mathbf{u}}\|_{2,\omega} \geq \\
&\geq \sup_{\eta \in \Upsilon(\omega)} \frac{\int (\mathbf{u} - \mathbf{v})\eta \, d\mathbf{x}}{\|\eta\|_{2,\omega}} - \sup_{\eta \in L^2(\omega)} \frac{\int (\hat{\mathbf{u}} - \mathbf{u})\eta \, d\mathbf{x}}{\|\eta\|_{2,\omega}} - \|\mathbf{u} - \hat{\mathbf{u}}\|_{2,\omega} = \\
&= |\mathbf{u} - \mathbf{v}|_{\Upsilon} - 2\|\mathbf{u} - \hat{\mathbf{u}}\|_{2,\omega}, \quad \forall \hat{\mathbf{u}} \in \Upsilon(\omega).
\end{aligned}$$

From here, it follows that

$$\|\mathbf{u} - \mathbf{v}\|_{2,\omega} \geq |\mathbf{u} - \mathbf{v}|_{\Upsilon} - 2\text{dist}(\mathbf{u}, \Upsilon(\omega)). \quad (132)$$

By (131) and (132), we conclude that the error arising if $\|\mathbf{u} - \mathbf{u}_h\|_{2,\Omega}$ is replaced by $|\mathbf{u} - \mathbf{u}_h|_{\Upsilon}$ depends on the regularity of \mathbf{u} and approximation properties of $\Upsilon(\omega)$.

Lecture 6

A POSTERIORI ESTIMATES FOR MIXED METHODS

Mixed approximations. A glance from the minimax theory

Consider our basic problem

$$\begin{aligned} \operatorname{div} \mathbf{A} \nabla \mathbf{u} + \mathbf{f} &= \mathbf{0} \quad \text{in } \Omega, \\ \mathbf{u} &= \mathbf{u}_0 \text{ on } \partial_1 \Omega, \\ \mathbf{A} \nabla \mathbf{u} \cdot \mathbf{n} &= \mathbf{F} \text{ on } \partial_2 \Omega, \end{aligned}$$

$$\mathbf{c}_1^2 |\xi|^2 \leq \mathbf{A}(\mathbf{x}) \xi \cdot \xi \leq \mathbf{c}_2^2 |\xi|^2 \quad \forall \xi \in \mathbb{R}^d, \text{ for a.e. } \mathbf{x} \in \Omega,$$

where $\mathbf{u}_0 \in \mathbf{H}^1(\Omega)$, $\mathbf{f} \in \mathbf{L}_2(\Omega)$, $\mathbf{F} \in \mathbf{L}_2(\partial_2 \Omega)$. Functional spaces

$$\begin{aligned} \mathbf{V} &:= \mathbf{H}^1(\Omega), \quad \mathbf{V}_0 := \{\mathbf{v} \in \mathbf{V} \mid \mathbf{v} = \mathbf{0} \text{ on } \partial_1 \Omega\}, \quad \widehat{\mathbf{V}} := \mathbf{L}_2(\Omega), \\ \mathbf{Q} &:= \mathbf{L}_2(\Omega; \mathbb{R}^d) \quad \widehat{\mathbf{Q}} := \mathbf{H}(\Omega; \operatorname{div}), \\ \widehat{\mathbf{Q}}^+ &:= \{\mathbf{y} \in \widehat{\mathbf{Q}} \mid \mathbf{y} \cdot \mathbf{n}|_{\partial_2 \Omega} \in \mathbf{L}_2(\partial_2 \Omega)\}. \end{aligned}$$

We recall that $\|\mathbf{q}\|_{\text{div}}$ is the norm in $H(\Omega; \text{div})$:

$$\|\mathbf{q}\|_{\text{div}} := (\|\mathbf{q}\|^2 + \|\text{div}\mathbf{q}\|^2)^{1/2} \quad \forall \mathbf{q} \in \mathbf{Q}$$

and

$$\|\mathbf{q}\| := \left(\int_{\Omega} \mathbf{A}\mathbf{q} \cdot \mathbf{q} \, dx \right)^{1/2}, \quad \mathbf{q} \in \mathbf{Q}$$

$$\|\mathbf{q}\|_* := \left(\int_{\Omega} \mathbf{A}^{-1}\mathbf{q} \cdot \mathbf{q} \, dx \right)^{1/2}$$

Note that,

$$\bar{c}_1^2 |\xi|^2 \leq \mathbf{A}^{-1}(\mathbf{x})\xi \cdot \xi \leq \bar{c}_2^2 |\xi|^2 \quad \forall \xi \in \mathbb{R}^d, \text{ for a.e. } \mathbf{x} \in \Omega$$

with $\bar{c}_1 = \mathbf{1}/c_2$, $\bar{c}_2 = \mathbf{1}/c_1$.

Generalized solution can be viewed as a saddle point of the Lagrangian

$$\mathbf{L}(\mathbf{v}, \mathbf{q}) := \int_{\Omega} \left(\nabla \mathbf{v} \cdot \mathbf{q} - \frac{1}{2} \mathbf{A}^{-1} \mathbf{q} \cdot \mathbf{q} \right) \mathbf{d}\mathbf{x} - \ell(\mathbf{v}),$$

where $\ell(\mathbf{v}) = \int_{\Omega} \mathbf{f}\mathbf{v} \mathbf{d}\mathbf{x} + \int_{\partial_2 \Omega} \mathbf{F}\mathbf{v} \mathbf{d}\mathbf{s}$.

The problem of finding $(\mathbf{u}, \mathbf{p}) \in \mathbf{V}_0 + \mathbf{u}_0 \times \mathbf{Q}$ such that

$$\mathbf{L}(\mathbf{u}, \mathbf{q}) \leq \mathbf{L}(\mathbf{u}, \mathbf{p}) \leq \mathbf{L}(\mathbf{v}, \mathbf{p}) \quad \forall \mathbf{q} \in \mathbf{Q}, \forall \mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0 \quad (133)$$

leads to is the so-called **Primal Mixed Formulation**

In this formulation the **solution** of a PDE is understood as **the pair of functions** $(\mathbf{u}, \mathbf{p}) \in (\mathbf{V}_0 + \mathbf{u}_0) \times \mathbf{Q}$ satisfying the relations

$$\int_{\Omega} (\mathbf{A}^{-1} \mathbf{p} - \nabla \mathbf{u}) \cdot \mathbf{q} \, dx = 0 \quad \forall \mathbf{q} \in \mathbf{Q}, \quad (134)$$

$$\int_{\Omega} \mathbf{p} \cdot \nabla \mathbf{w} \, dx - \ell(\mathbf{w}) = 0 \quad \forall \mathbf{w} \in \mathbf{V}_0. \quad (135)$$

In the PMM,

$\mathbf{p} = \mathbf{A} \nabla \mathbf{u}$, is satisfied in $\mathbf{L}_2(\Omega)$ – sense

$\operatorname{div} \mathbf{p} + \mathbf{f} = \mathbf{0}$ in Ω and $\mathbf{p} \cdot \mathbf{n} = \mathbf{F}$ on $\partial_2 \Omega$ are satisfied in a weak sense.

As we have seen in previous lectures \mathbf{L} generates two functionals

$$\mathbf{J}(\mathbf{v}) := \sup_{\mathbf{q} \in \mathbf{Q}} \mathbf{L}(\mathbf{v}, \mathbf{q}) = \frac{1}{2} \|\nabla \mathbf{v}\|^2 - \ell(\mathbf{v})$$

and

$$\mathbf{I}^*(\mathbf{q}) := -\frac{1}{2} \|\mathbf{q}\|_*^2 - \ell(\mathbf{u}_0) + \int_{\Omega} \nabla \mathbf{u}_0 \cdot \mathbf{q} \, \mathbf{d}\mathbf{x}.$$

Also, we know that

$$\inf_{\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0} \mathbf{J}(\mathbf{v}) := \inf \mathcal{P} = \mathbf{L}(\mathbf{u}, \mathbf{p}) = \sup \mathcal{P}^* := \sup_{\mathbf{q} \in \mathbf{Q}_{\ell}} \mathbf{I}^*(\mathbf{q}), \quad (136)$$

where $\mathbf{Q}_{\ell} := \{\mathbf{q} \in \mathbf{Q} \mid \int_{\Omega} \mathbf{q} \cdot \nabla \mathbf{w} \, \mathbf{d}\mathbf{x} = \ell(\mathbf{w}) \quad \forall \mathbf{w} \in \mathbf{V}_0\}$.

Primal Mixed Method (PMM)

Let $\mathbf{Q}_h \subset \mathbf{Q}$ and $\mathbf{V}_{0h} \subset \mathbf{V}_0$ are subspaces constructed by FE approximation, then a discrete analog of (134)–(135) is the

Primal Mixed Finite Element Method

See, e.g., F. Brezzi and M. Fortin, *Mixed and Hybrid Finite Element Methods*. Springer-Verlag, New York, 1991.

D. Braess. *Finite elements*. Cambridge University Press, Cambridge, 1997.

J. E. Roberts and J.-M. Thomas, *Mixed and Hybrid Methods*. In Handbook of Numerical Analysis, II, eds. P. G. Ciarlet and J.-L. Lions, North-Holland, Amsterdam, pp. 523–639, 1991.

In PMM, we need to find a pair of functions $(\mathbf{u}_h, \mathbf{p}_h) \in (\mathbf{V}_{0h} + \mathbf{u}_0) \times \mathbf{Q}_h$ such that

$$\int_{\Omega} (\mathbf{A}^{-1} \mathbf{p}_h - \nabla \mathbf{u}_h) \cdot \mathbf{q}_h \, d\mathbf{x} = 0 \quad \forall \mathbf{q}_h \in \mathbf{Q}_h, \quad (137)$$

$$\int_{\Omega} \mathbf{p}_h \cdot \nabla \mathbf{w}_h \, d\mathbf{x} - \ell(\mathbf{w}_h) = 0 \quad \forall \mathbf{w}_h \in \mathbf{V}_{0h}. \quad (138)$$

In this formulation, \mathbf{u}_h can be constructed by means of the Courant-type elements and \mathbf{p}_h by piecewise constant functions.

Dual Mixed Method (DMM)

Another mixed formulation arises if we represent \mathbf{L} in a somewhat different form. First, we introduce the functional $\mathbf{g} : (\mathbf{V}_0 + \mathbf{u}_0) \times \widehat{\mathbf{Q}} \rightarrow \mathbb{R}$ by the relation

$$\mathbf{g}(\mathbf{v}, \mathbf{q}) := \int_{\Omega} (\nabla \mathbf{v} \cdot \mathbf{q} + \mathbf{v}(\operatorname{div} \mathbf{q})) \, dx.$$

We have

$$\begin{aligned} \mathbf{L}(\mathbf{v}, \mathbf{q}) &= \int_{\Omega} \left(\nabla \mathbf{v} \cdot \mathbf{q} - \frac{1}{2} \mathbf{A}^{-1} \mathbf{q} \cdot \mathbf{q} \right) \, dx - \ell(\mathbf{v}) = \\ &= \mathbf{g}(\mathbf{v}, \mathbf{q}) - \int_{\Omega} \mathbf{v}(\operatorname{div} \mathbf{q}) \, dx - \frac{1}{2} \|\mathbf{q}\|_*^2 - \ell(\mathbf{v}). \end{aligned}$$

Introduce the set

$$\widehat{\mathbf{Q}}_{\mathbf{F}} := \{\mathbf{q} \in \widehat{\mathbf{Q}} \mid \mathbf{g}(\mathbf{w}, \mathbf{q}) = \int_{\partial_2 \Omega} \mathbf{F} \mathbf{w} \, ds \quad \forall \mathbf{w} \in \mathbf{V}_0\}.$$

Note that for $\mathbf{q} \in \widehat{\mathbf{Q}}_{\mathbf{F}}$ we have

$$\begin{aligned} \mathbf{g}(\mathbf{v}, \mathbf{q}) &= \mathbf{g}(\mathbf{w} + \mathbf{u}_0, \mathbf{q}) = \mathbf{g}(\mathbf{w}, \mathbf{q}) + \mathbf{g}(\mathbf{u}_0, \mathbf{q}) = \\ &= \int_{\partial_2 \Omega} \mathbf{F} \mathbf{w} \, ds + \mathbf{g}(\mathbf{u}_0, \mathbf{q}) \quad \forall \mathbf{w} \in \mathbf{V}_0. \end{aligned}$$

Therefore, if the variable \mathbf{q} is taken not from \mathbf{Q} but from the narrower set $\widehat{\mathbf{Q}}_{\mathbf{F}}$, then the Lagrangian can be written as

$$\begin{aligned} \widehat{\mathbf{L}}(\mathbf{v}, \mathbf{q}) &:= \\ &= -\frac{1}{2} \|\mathbf{q}\|_*^2 - \int_{\Omega} \mathbf{v}(\operatorname{div} \mathbf{q}) \, dx - \int_{\Omega} \mathbf{f} \mathbf{v} \, dx - \int_{\partial_2 \Omega} \mathbf{F} \mathbf{u}_0 \, ds + \mathbf{g}(\mathbf{u}_0, \mathbf{q}). \end{aligned}$$

We observe Note the new Lagrangian $\widehat{\mathbf{L}}$
 is defined on a wider set of primal functions $\mathbf{v} \in \widehat{\mathbf{V}}$, but uses a narrower set
 $\widehat{\mathbf{Q}}_F$ for the fluxes.

$$\begin{aligned} \text{The problem of finding } (\widehat{\mathbf{u}}, \widehat{\mathbf{p}}) \in \widehat{\mathbf{V}} \times \widehat{\mathbf{Q}}_F \text{ such that} \\ \widehat{\mathbf{L}}(\widehat{\mathbf{u}}, \widehat{\mathbf{q}}) \leq \widehat{\mathbf{L}}(\widehat{\mathbf{u}}, \widehat{\mathbf{p}}) \leq \widehat{\mathbf{L}}(\widehat{\mathbf{v}}, \widehat{\mathbf{p}}) \quad \forall \widehat{\mathbf{q}} \in \widehat{\mathbf{Q}}_F, \forall \widehat{\mathbf{v}} \in \widehat{\mathbf{V}} \end{aligned} \quad (139)$$

lead to is the so-called

Dual Mixed Formulation

of the problem in question (see, e.g., F. Brezzi and M. Fortin).

From (139) we obtain the necessary conditions for the dual mixed formulation. Since

$$\widehat{\mathbf{L}}(\widehat{\mathbf{u}}, \widehat{\mathbf{q}}) \leq \widehat{\mathbf{L}}(\widehat{\mathbf{u}}, \widehat{\mathbf{p}}) \quad \forall \widehat{\mathbf{q}} \in \widehat{\mathbf{Q}}_F,$$

we have

$$\begin{aligned} -\frac{1}{2} \|\widehat{\mathbf{p}} + \lambda \boldsymbol{\eta}\|_*^2 - \int_{\Omega} \widehat{\mathbf{u}}(\operatorname{div}(\widehat{\mathbf{p}} + \lambda \boldsymbol{\eta}) - \widehat{\mathbf{f}}) \, dx - \int_{\partial_2 \Omega} \mathbf{F} \mathbf{u}_0 \, ds + \mathbf{g}(\mathbf{u}_0, \widehat{\mathbf{p}} + \lambda \boldsymbol{\eta}) \leq \\ -\frac{1}{2} \|\widehat{\mathbf{p}}\|_*^2 - \int_{\Omega} \widehat{\mathbf{u}}(\operatorname{div} \widehat{\mathbf{p}}) \, dx - \int_{\Omega} \widehat{\mathbf{f}} \, dx - \int_{\partial_2 \Omega} \mathbf{F} \mathbf{u}_0 \, ds + \mathbf{g}(\mathbf{u}_0, \widehat{\mathbf{p}}), \end{aligned}$$

where λ is a real number and $\boldsymbol{\eta}$ is a function in $\widehat{\mathbf{Q}}_0 := \widehat{\mathbf{Q}}_F$ with $\mathbf{F} = \mathbf{0}$. Now, arrive at the relation

$$-\lambda \int_{\Omega} (\mathbf{A}^{-1} \widehat{\mathbf{p}} \cdot \boldsymbol{\eta} + \widehat{\mathbf{u}}(\operatorname{div} \boldsymbol{\eta})) \, dx + \lambda \mathbf{g}(\mathbf{u}_0, \boldsymbol{\eta}) \leq \frac{\lambda^2}{2} \int_{\Omega} \mathbf{A}^{-1} \boldsymbol{\eta} \cdot \boldsymbol{\eta} \, dx.$$

Rewrite it as

$$\int_{\Omega} (\mathbf{A}^{-1} \hat{\mathbf{p}} \cdot \boldsymbol{\eta} + \hat{\mathbf{u}}(\operatorname{div} \boldsymbol{\eta})) \, \mathrm{d}\mathbf{x} - \mathbf{g}(\mathbf{u}_0, \boldsymbol{\eta}) \geq \frac{\lambda}{2} \int_{\Omega} \mathbf{A}^{-1} \boldsymbol{\eta} \cdot \boldsymbol{\eta} \, \mathrm{d}\mathbf{x}.$$

Since $\lambda > 0$ can be taken arbitrarily small, the latter relation may hold only if

$$\int_{\Omega} (\mathbf{A}^{-1} \hat{\mathbf{p}} \cdot \boldsymbol{\eta} + \hat{\mathbf{u}} \operatorname{div} \boldsymbol{\eta}) \, \mathrm{d}\mathbf{x} - \mathbf{g}(\mathbf{u}_0, \boldsymbol{\eta}) \geq 0.$$

But $\boldsymbol{\eta}$ is an arbitrary element of a linear manifold $\widehat{\mathbf{Q}}_0$, so that $+\boldsymbol{\eta}$ can be replaced by $-\boldsymbol{\eta}$ what leads to the conclusion that

$$\int_{\Omega} (\mathbf{A}^{-1} \hat{\mathbf{p}} \cdot \boldsymbol{\eta} + \hat{\mathbf{u}} \operatorname{div} \boldsymbol{\eta}) \, \mathrm{d}\mathbf{x} - \mathbf{g}(\mathbf{u}_0, \boldsymbol{\eta}) = 0 \quad \forall \boldsymbol{\eta} \in \widehat{\mathbf{Q}}_0.$$

From

$$\widehat{\mathbf{L}}(\widehat{\mathbf{u}}, \widehat{\mathbf{p}}) \leq \widehat{\mathbf{L}}(\widehat{\mathbf{u}} + \widehat{\mathbf{v}}, \widehat{\mathbf{p}}) \quad \forall \widehat{\mathbf{v}} \in \widehat{\mathbf{V}} := \mathbf{L}^2(\Omega)$$

we observe that the terms of $\widehat{\mathbf{L}}$ linear with respect to the "pressure" must vanish. Namely, we obtain

$$\int_{\Omega} (\widehat{\mathbf{v}} \operatorname{div} \widehat{\mathbf{p}} + \widehat{\mathbf{f}} \widehat{\mathbf{v}}) \, \mathbf{d}\mathbf{x} = 0$$

Thus, we arrive at the system

$$\int_{\Omega} (\mathbf{A}^{-1} \widehat{\mathbf{p}} \cdot \widehat{\mathbf{q}} + (\operatorname{div} \widehat{\mathbf{q}}) \widehat{\mathbf{u}}) \, \mathbf{d}\mathbf{x} = \mathbf{g}(\mathbf{u}_0, \widehat{\mathbf{q}}) \quad \forall \widehat{\mathbf{q}} \in \widehat{\mathbf{Q}}_0, \quad (140)$$

$$\int_{\Omega} (\operatorname{div} \widehat{\mathbf{p}} + \widehat{\mathbf{f}}) \widehat{\mathbf{v}} \, \mathbf{d}\mathbf{x} = 0 \quad \forall \widehat{\mathbf{v}} \in \widehat{\mathbf{V}}. \quad (141)$$

We observe that now the condition

$$\mathbf{div} \hat{\mathbf{p}} + \mathbf{f} = \mathbf{0}$$

is satisfied in a "strong" (\mathbf{L}_2) sense, the Neumann type boundary condition is viewed as the essential boundary condition, and the relation

$$\hat{\mathbf{p}} = \mathbf{A} \nabla \hat{\mathbf{u}}$$

and the Dirichlet type boundary condition are satisfied in a weak sense. These properties of the DMM lead to that the respective finite dimensional formulations are better adapted to the satisfaction of the equilibrium type relations for the fluxes. This fact is important in many applications where a sharp satisfaction of the equilibrium relations is required.

The Lagrangian $\widehat{\mathbf{L}}$ also generates two functionals

$$\widehat{\mathbf{J}}(\widehat{\mathbf{v}}) := \sup_{\widehat{\mathbf{q}} \in \widehat{\mathbf{Q}}_F} \widehat{\mathbf{L}}(\widehat{\mathbf{v}}, \widehat{\mathbf{q}}) \quad \text{and} \quad \widehat{\mathbf{I}}^*(\widehat{\mathbf{q}}) := \inf_{\widehat{\mathbf{v}} \in \widehat{\mathbf{V}}} \widehat{\mathbf{L}}(\widehat{\mathbf{v}}, \widehat{\mathbf{q}}).$$

The two corresponding variational problems are

$$\inf_{\widehat{\mathbf{v}} \in \widehat{\mathbf{V}}} \widehat{\mathbf{J}}(\widehat{\mathbf{v}}) \quad \text{and} \quad \sup_{\widehat{\mathbf{q}} \in \widehat{\mathbf{Q}}_F} \widehat{\mathbf{I}}^*(\widehat{\mathbf{q}}).$$

They are called Problems $\widehat{\mathcal{P}}$ and $\widehat{\mathcal{P}}^*$, respectively. Note that the functional $\widehat{\mathbf{J}}$ (unlike \mathbf{J}) has no simple explicit form. However, we can prove the solvability of Problem $\widehat{\mathcal{P}}$ by the following Lemma.

Lemma

For any $\widehat{\mathbf{v}} \in \widehat{\mathbf{V}}$ and $\mathbf{F} \in \mathbf{L}_2(\partial_2\Omega)$ there exists $\mathbf{p}^{\mathbf{v}} \in \widehat{\mathbf{Q}}_{\mathbf{F}}$ such that

$$\mathbf{div}^{\mathbf{v}} + \widehat{\mathbf{v}} = \mathbf{0} \quad \text{in } \Omega, \quad (142)$$

$$\|\mathbf{p}^{\mathbf{v}}\|_* \leq \mathbf{C}_{\Omega} (\|\widehat{\mathbf{v}}\| + \|\mathbf{F}\|_{\partial_2\Omega}). \quad (143)$$

Proof. We know that the boundary-value problem

$$\begin{aligned} \mathbf{div} \mathbf{A} \nabla \mathbf{u}^{\mathbf{v}} + \widehat{\mathbf{v}} &= \mathbf{0} && \text{in } \Omega, \\ \mathbf{u}^{\mathbf{v}} &= \mathbf{0} && \text{on } \partial_1\Omega, \\ \mathbf{A} \nabla \mathbf{u}^{\mathbf{v}} \cdot \mathbf{n} &= \mathbf{F} && \text{on } \partial_2\Omega \end{aligned}$$

possesses the unique solution $\mathbf{u}^{\mathbf{v}} \in \mathbf{V}_0$.

For it and the energy estimate

$$\| \nabla \mathbf{u}^v \| \leq \mathbf{C}_\Omega (\| \hat{\mathbf{v}} \| + \| \mathbf{F} \|_{\partial_2 \Omega})$$

holds. Let $\mathbf{p}^v := \mathbf{A} \nabla \mathbf{u}^v$. We have

$$\operatorname{div} \mathbf{p}^v + \hat{\mathbf{v}} = \mathbf{0}.$$

Obviously, $\mathbf{p}^v \in \hat{\mathbf{Q}}_F$ and, since

$$\| \mathbf{p}^v \|_*^2 = \int_{\Omega} \mathbf{A}^{-1} (\mathbf{A} \nabla \mathbf{u}^v) \cdot (\mathbf{A} \nabla \mathbf{u}^v) \, dx = \| \nabla \mathbf{u}^v \|^2,$$

we find that (143) also holds.

□

By the Lemma we can easily prove the coercivity of $\widehat{\mathbf{J}}$ on $\widehat{\mathbf{V}}$. Indeed,

$$\begin{aligned} \widehat{\mathbf{J}}(\widehat{\mathbf{v}}) &\geq \widehat{\mathbf{L}}(\widehat{\mathbf{v}}, \alpha \mathbf{p}^{\mathbf{v}}) = \\ &= -\frac{1}{2} \|\alpha \mathbf{p}^{\mathbf{v}}\|_*^2 - \alpha \int_{\Omega} \widehat{\mathbf{v}}(\operatorname{div} \mathbf{p}^{\mathbf{v}}) \, dx - \int_{\Omega} \mathbf{f} \widehat{\mathbf{v}} \, dx - \int_{\partial_2 \Omega} \mathbf{F} \mathbf{u}_0 \, ds + \mathbf{g}(\mathbf{u}_0, \alpha \mathbf{p}^{\mathbf{v}}) = \\ &= -\frac{1}{2} \alpha^2 \|\mathbf{p}^{\mathbf{v}}\|_*^2 + \alpha \|\widehat{\mathbf{v}}\|^2 - \|\mathbf{f}\| \|\widehat{\mathbf{v}}\| + \mathbf{g}(\mathbf{u}_0, \alpha \mathbf{p}^{\mathbf{v}}) - \int_{\partial_2 \Omega} \mathbf{F} \mathbf{u}_0 \, ds. \end{aligned}$$

Here $|\mathbf{g}(\mathbf{u}_0, \alpha \mathbf{p}^{\mathbf{v}})| \leq \alpha \|\mathbf{p}^{\mathbf{v}}\|_{\operatorname{div}} \|\mathbf{u}_0\|_{1,2,\Omega}$ and

$$\begin{aligned} \|\mathbf{p}^{\mathbf{v}}\|_{\operatorname{div}}^2 &= \|\mathbf{p}^{\mathbf{v}}\|^2 + \|\operatorname{div} \mathbf{p}^{\mathbf{v}}\|^2 \leq \frac{1}{\underline{c}_1} \|\mathbf{p}^{\mathbf{v}}\|_*^2 + \|\widehat{\mathbf{v}}\|^2 \leq \\ &\leq \frac{1}{\underline{c}_1} \mathbf{C}_{\Omega}^2 (\|\widehat{\mathbf{v}}\| + \|\mathbf{F}\|_{\partial_2 \Omega})^2 + \|\widehat{\mathbf{v}}\|^2. \end{aligned}$$

Therefore

$$\hat{J}(\hat{\mathbf{v}}) \geq -\frac{1}{2}\alpha^2 \mathbf{C}_\Omega^2 \|\hat{\mathbf{v}}\|^2 + \alpha \|\hat{\mathbf{v}}\|^2 + \Theta(\|\hat{\mathbf{v}}\|) + \Theta_0,$$

where $\Theta(\|\hat{\mathbf{v}}\|)$ contains the terms linear with respect to $\|\hat{\mathbf{v}}\|$ and Θ_0 does not depend on $\hat{\mathbf{v}}$. Take $\alpha = \mathbf{1}/\mathbf{C}_\Omega^2$. Then

$$\hat{J}(\hat{\mathbf{v}}) \geq \frac{1}{2\mathbf{C}_\Omega^2} \|\hat{\mathbf{v}}\|^2 + \Theta(\|\hat{\mathbf{v}}\|) + \Theta_0 \longrightarrow +\infty \text{ as } \|\hat{\mathbf{v}}\| \rightarrow \infty.$$

It is not difficult to prove that the functional \hat{J} is convex and lower semicontinuous. Therefore, Problem \hat{P} has a solution $\hat{\mathbf{u}}$.

Inf-Sup condition for the dual mixed formulation

Corollary

Lemma implies the *inf-sup* condition

$$\inf_{\substack{\phi \in L^2(\Omega) \\ \psi \in L^2(\partial_2\Omega)}} \sup_{\mathbf{q} \in \widehat{\mathbf{Q}}_F} \frac{\int_{\Omega} \phi \operatorname{div} \mathbf{q} \, dx + \int_{\partial_2\Omega} \psi \mathbf{q} \cdot \mathbf{n} \, ds}{\|\mathbf{q}\|_{\operatorname{div}} (\|\phi\|^2 + \|\psi\|_{\partial_2\Omega}^2)^{1/2}} \geq \mathbf{C}_0 > \mathbf{0}.$$

The Dual Problem with respect to the Lagrangian $\widehat{\mathbf{L}}$

Let us now construct the dual functional $\widehat{\mathbf{I}}^*$. It is easy to see that

$$\begin{aligned} \widehat{\mathbf{I}}^*(\widehat{\mathbf{q}}) &= \inf_{\widehat{\mathbf{v}}} \widehat{\mathbf{L}}(\widehat{\mathbf{v}}, \widehat{\mathbf{q}}) = \\ &= \inf_{\widehat{\mathbf{v}}} \left\{ -\frac{1}{2} \|\widehat{\mathbf{q}}\|_*^2 - \int_{\Omega} \mathbf{v}(\operatorname{div}\widehat{\mathbf{q}}) dx - \int_{\Omega} \mathbf{f} \mathbf{v} dx - \int_{\partial_2\Omega} \mathbf{F} \mathbf{u}_0 ds + \mathbf{g}(\mathbf{u}_0, \widehat{\mathbf{q}}) \right\} = \\ &= -\frac{1}{2} \|\widehat{\mathbf{q}}\|_*^2 + \mathbf{g}(\mathbf{u}_0, \widehat{\mathbf{q}}) - \int_{\partial_2\Omega} \mathbf{F} \mathbf{u}_0 ds \end{aligned}$$

provided that $\operatorname{div}\widehat{\mathbf{q}} + \mathbf{f} = \mathbf{0}$ (in the \mathbf{L}_2 -sense). In all other cases $\widehat{\mathbf{I}}^*(\widehat{\mathbf{q}}) = -\infty$.

Recalling that $\mathbf{div} \hat{\mathbf{q}} = -\mathbf{f}$ (in $L_2(\Omega)$ -sense), we find that the dual functional for such a case has the form

$$\begin{aligned} \hat{\mathbf{I}}^*(\mathbf{q}) &= -\frac{1}{2} \|\hat{\mathbf{q}}\|_*^2 + \int_{\Omega} (\nabla \mathbf{u}_0 \cdot \hat{\mathbf{q}} - \mathbf{f} \mathbf{u}_0) \, dx - \int_{\partial_2 \Omega} \mathbf{F} \mathbf{u}_0 \, ds \\ &= \int_{\Omega} \nabla \mathbf{u}_0 \cdot \hat{\mathbf{q}} \, dx - \frac{1}{2} \|\hat{\mathbf{q}}\|_*^2 - \ell(\mathbf{u}_0), \end{aligned}$$

Since $\hat{\mathbf{q}} \in \hat{\mathbf{Q}}_{\mathbf{F}}$, we have

$$\int_{\Omega} \nabla \mathbf{w} \cdot \hat{\mathbf{q}} \, dx = - \int_{\Omega} (\mathbf{div} \hat{\mathbf{q}}) \mathbf{w} \, dx + \int_{\partial_2 \Omega} \mathbf{F} \mathbf{w} \, ds \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

we see that $\hat{\mathbf{q}}$ satisfies the relation

$$\int_{\Omega} \nabla \mathbf{w} \cdot \hat{\mathbf{q}} \, dx = \ell(\mathbf{w}) \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

In other cases, $\hat{\mathbf{I}}^*(\hat{\mathbf{q}}) = -\infty$.

Thus, Problems \mathcal{P}^* and $\widehat{\mathcal{P}}^*$ coincide and are reduced to the maximization of \mathbf{I}^* on the set Q_ℓ . This means that

$$\sup \mathcal{P}^* = \sup \widehat{\mathcal{P}}^*.$$

Since the saddle point of $\widehat{\mathbf{L}}$ exists, we have

$$\widehat{\mathbf{L}}(\widehat{\mathbf{u}}, \widehat{\mathbf{p}}) = \inf \widehat{\mathcal{P}} = \sup \widehat{\mathcal{P}}^*,$$

but

$$\sup \widehat{\mathcal{P}}^* = \sup \mathcal{P}^* = \inf \mathcal{P}.$$

Thus, we infer that

$$\inf \widehat{\mathcal{P}} = \inf \mathcal{P}.$$

Thus, we conclude that $\mathbf{u} \in \mathbf{V}_0 + \mathbf{u}_0$ (minimizer of \mathcal{P}) also minimizes $\widehat{\mathbf{J}}$ on $\widehat{\mathbf{V}}$. Analogously, if $\mathbf{p} \in \mathbf{Q}_\ell$ is the maximizer of Problem \mathcal{P}^* , then

$$\int_{\Omega} \nabla \mathbf{w} \cdot \mathbf{p} \, \mathbf{d}\mathbf{x} = \int_{\Omega} \mathbf{f}\mathbf{w} \, \mathbf{d}\mathbf{x} + \int_{\partial_2\Omega} \mathbf{F}\mathbf{w} \, \mathbf{d}\mathbf{s} \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

From here we see that $\mathbf{divp} + \mathbf{f} = \mathbf{0}$ a.e. in Ω and, hence,

$$\int_{\Omega} (\nabla \mathbf{w} \cdot \mathbf{p} + (\mathbf{divp})\mathbf{w}) \, \mathbf{d}\mathbf{x} = \int_{\partial_2\Omega} \mathbf{F}\mathbf{w} \, \mathbf{d}\mathbf{s} \quad \forall \mathbf{w} \in \mathbf{V}_0,$$

that is $\mathbf{p} \in \widehat{\mathbf{Q}}_F$. Thus, \mathbf{p} is also the maximizer of Problem $\widehat{\mathcal{P}}^*$.

The reverse statement that the solutions of $\widehat{\mathcal{P}}$, $\widehat{\mathcal{P}}^*$ are also the solutions of \mathcal{P} , \mathcal{P}^* is not difficult to prove as well.

Hence, both mixed formulations have the same solution (u, p) which is in fact the generalized solution of our problem.

Finite dimensional formulations

Let

$$\widehat{\mathbf{V}}_h \subset \widehat{\mathbf{V}}, \quad \widehat{\mathbf{Q}}_{0h} \subset \widehat{\mathbf{Q}}_0, \quad \widehat{\mathbf{Q}}_{Fh} \subset \widehat{\mathbf{Q}}_F$$

A discrete analog of the dual mixed formulation is: Find $(\widehat{\mathbf{u}}_h, \widehat{\mathbf{p}}_h) \in \widehat{\mathbf{V}}_h \times \widehat{\mathbf{Q}}_{Fh}$ such that

$$\int_{\Omega} (\mathbf{A}^{-1} \widehat{\mathbf{p}}_h \cdot \widehat{\mathbf{q}}_h + \widehat{\mathbf{u}}_h \operatorname{div} \widehat{\mathbf{q}}_h) \, dx = \mathbf{g}(\mathbf{u}_0, \widehat{\mathbf{q}}_h) \quad \forall \widehat{\mathbf{q}}_h \in \widehat{\mathbf{Q}}_{0h}, \quad (144)$$

$$\int_{\Omega} (\operatorname{div} \widehat{\mathbf{p}}_h + \mathbf{f}) \widehat{\mathbf{v}}_h \, dx = 0 \quad \forall \widehat{\mathbf{v}}_h \in \widehat{\mathbf{V}}_h. \quad (145)$$

Error analysis for DMM

First we will obtain a priori error estimates for the dual mixed method and after that we will derive computable upper bounds for the quantities

$$\| \nabla(\mathbf{u} - \mathbf{u}_h) \|, \| \mathbf{p} - \mathbf{p}_h \|_*, \| \mathbf{p} - \hat{\mathbf{p}}_h \|_{\text{div}} .$$

A priori error estimates for DMM

Below we will show a simple way of the derivation of projection type error estimates for the dual mixed method. By combining them with standard interpolation results, one can obtain known rate convergence estimates. A detailed exposition of this subject can be found in the above cited books. Here, we present a simplified version, which, however contains the main ideas of the a priori error analysis for the dual mixed approximations.

For the sake of simplicity we will consider the case of uniform Dirichlet boundary conditions and a constant matrix \mathbf{A} . In this case, the basic system is as follows

$$\int_{\Omega} (\mathbf{A}^{-1} \hat{\mathbf{p}} \cdot \hat{\mathbf{q}} + (\operatorname{div} \hat{\mathbf{q}}) \hat{u}) \, dx = 0 \quad \forall \hat{\mathbf{q}} \in \hat{\mathbf{Q}}_0,$$

$$\int_{\Omega} (\operatorname{div} \hat{\mathbf{p}} + \mathbf{f}) \hat{v} \, dx = 0 \quad \forall \hat{v} \in \hat{\mathbf{V}}.$$

Since there is no Neumann part of the boundary, $\hat{\mathbf{Q}}_F$ and $\hat{\mathbf{Q}}_0$ coincides with $\hat{\mathbf{Q}} := \mathbf{H}(\Omega, \operatorname{div})$.

In the considered, case the system of DMM is as follows

$$\int_{\Omega} \left(\mathbf{A}^{-1} \hat{\mathbf{p}}_h \cdot \hat{\mathbf{q}}_h + \hat{\mathbf{u}}_h \operatorname{div} \hat{\mathbf{q}}_h \right) dx = 0 \quad \forall \hat{\mathbf{q}}_h \in \hat{\mathbf{Q}}_h,$$

$$\int_{\Omega} (\operatorname{div} \hat{\mathbf{p}}_h + \mathbf{f}) \hat{\mathbf{v}}_h dx = 0 \quad \forall \hat{\mathbf{v}}_h \in \hat{\mathbf{V}}_h.$$

Assumptions.

- (a) \mathcal{T}_h is a regular triangulation of a polygonal domain Ω .
- (b) $\hat{\mathbf{V}}_h = \{\mathbf{v}_h \in \mathbf{L}^2 \mid \mathbf{v}_h \in \mathbf{P}^0(\mathbf{T}) \quad \forall \mathbf{T} \in \mathcal{T}_h\}$.
- (c) $\hat{\mathbf{Q}}_h = \{\mathbf{q}_h \in \mathbf{H}(\Omega, \operatorname{div}) \mid \mathbf{q}_h \in \mathbf{RT}^0(\mathbf{T}) \quad \forall \mathbf{T} \in \mathcal{T}_h\}$.
- (d) $\mathbf{f} \in \mathbf{P}^0(\mathbf{T}), \quad \forall \mathbf{T} \in \mathcal{T}_h$

Note that under the assumptions made

$$\mathbf{div} \mathbf{p}_h + \mathbf{f} = \mathbf{0} \quad \text{on any } \mathbf{T}.$$

Indeed, this fact directly follows from the relation

$$\int_{\Omega} (\mathbf{div} \hat{\mathbf{p}}_h + \mathbf{f}) \hat{\mathbf{v}}_h \, d\mathbf{x} = \mathbf{0} \quad \forall \hat{\mathbf{v}}_h \in \hat{\mathbf{V}}_h.$$

Therefore $\mathbf{p}_h \in \mathbf{Q}_f$.

Compatibility and stability conditions

In order to provide the stability of the discrete DM formulation we need additional assumptions.

We assume that a pair of finite dimensional spaces $\widehat{\mathbf{V}}_h$, $\widehat{\mathbf{Q}}_h$ satisfies the following condition:

For any $\mathbf{v}_h \in \widehat{\mathbf{V}}_h$ exists $\mathbf{q}_h^v \in \widehat{\mathbf{Q}}_h$ such that

$$\mathbf{div} \mathbf{q}_h^v = \mathbf{v}_h \quad (\text{compatibility}), \quad (146)$$

$$\|\mathbf{q}_h^v\| \leq \mathbf{C} \|\mathbf{v}_h\| \quad (\text{stability}). \quad (147)$$

We will show that the above two conditions are the sufficient conditions for proving that discrete DM problem is

- (a) correct (e.g., has a solution),
- (b) stable,
- (c) has a projection type error estimate.

Discrete Inf-Sup condition

From (146) and (147), it follows that

$$\inf_{\mathbf{v}_h \in \widehat{\mathbf{V}}_h} \sup_{\mathbf{q}_h \in \widehat{\mathbf{Q}}_h} \frac{\int_{\Omega} \mathbf{v}_h \operatorname{div} \mathbf{q}_h \, dx}{\|\mathbf{v}_h\| \|\mathbf{q}_h\|_{\operatorname{div}}} \geq \mathbf{C} > \mathbf{0}$$

Indeed,

$$\sup_{\mathbf{q}_h \in \widehat{\mathbf{Q}}_h} \frac{\int_{\Omega} \mathbf{v}_h \operatorname{div} \mathbf{q}_h \, dx}{\|\mathbf{v}_h\| \|\mathbf{q}_h\|_{\operatorname{div}}} \geq \frac{\int_{\Omega} \mathbf{v}_h \operatorname{div} \mathbf{q}_h^{\vee} \, dx}{\|\mathbf{v}_h\| \|\mathbf{q}_h^{\vee}\|_{\operatorname{div}}} = \frac{\|\mathbf{v}_h\|}{\|\mathbf{q}_h\|_{\operatorname{div}}} \geq \frac{\mathbf{1}}{\sqrt{\mathbf{1} + \mathbf{C}^2}}.$$

Now, we refer to known results on the solvability of DMM, that can be summarized as follows: **if the triangulations are "regular" and the discrete Inf-Sup condition holds, then the discrete formulation has a unique solution.**

Projection type estimate for the dual problem

Since \mathbf{p} is a maximizer, i.e.,

$$-\frac{1}{2} \|\mathbf{q}\|_*^2 \leq -\frac{1}{2} \|\mathbf{p}\|_*^2 \quad \forall \mathbf{q} \in \mathbf{Q}_f,$$

we find that

$$\int_{\Omega} \mathbf{A}^{-1} \mathbf{p} \cdot \mathbf{q} \, dx = 0 \quad \forall \mathbf{q} \in \mathbf{Q}_0,$$

where \mathbf{Q}_0 is the space of solenoidal functions. Therefore, for any $\mathbf{q} \in \mathbf{Q}_f$,

$$\begin{aligned} \frac{1}{2} \|\mathbf{q} - \mathbf{p}\|_*^2 &= \frac{1}{2} \|\mathbf{q}\|_*^2 - \frac{1}{2} \|\mathbf{p}\|_*^2 + \int_{\Omega} \mathbf{A}^{-1} \mathbf{p} \cdot (\mathbf{p} - \mathbf{q}) \, dx = \\ &= \frac{1}{2} \|\mathbf{q}\|_*^2 - \frac{1}{2} \|\mathbf{p}\|_*^2. \end{aligned}$$

Let $\mathbf{Q}_{fh} = \mathbf{Q}_f \cap \widehat{\mathbf{Q}}_h$. Note that $\mathbf{p}_h \in \mathbf{Q}_{fh}$ is also the maximizer of $-\frac{1}{2} \|\mathbf{q}_{fh}\|_*^2$ on \mathbf{Q}_{fh} , so that

$$\begin{aligned} \frac{1}{2} \|\mathbf{p}_h - \mathbf{p}\|_*^2 &= \frac{1}{2} \|\mathbf{p}_h\|_*^2 - \frac{1}{2} \|\mathbf{p}\|_*^2 \leq \frac{1}{2} \|\mathbf{q}_{fh}\|_*^2 - \frac{1}{2} \|\mathbf{p}\|_*^2 = \\ &= \frac{1}{2} \|\mathbf{q}_{fh} - \mathbf{p}\|_*^2 \quad \forall \mathbf{q}_{fh} \in \mathbf{Q}_{fh}. \end{aligned}$$

Thus, we arrive at the first projection estimate

$$\|\mathbf{p} - \mathbf{p}_h\|_* \leq \inf_{\mathbf{q}_{fh} \in \mathbf{Q}_{fh}} \|\mathbf{p} - \mathbf{q}_{fh}\|_* . \quad (148)$$

However, this projection error estimate has an obvious drawback. It is applicable only for a very narrow class of approximations: conforming (internal) approximations of the set \mathbf{Q}_f .

To obtain an estimate for a wider class, we first derive one auxiliary result.

A Modified DM problem

Take $\tilde{\mathbf{f}} = \mathbf{div}(\hat{\mathbf{q}}_h - \mathbf{p})$ where $\hat{\mathbf{q}}_h \in \hat{\mathbf{Q}}_h$ and solve the **modified DM problem**

$$\int_{\Omega} \left(\mathbf{A}^{-1} \hat{\mathbf{p}}_h^f \cdot \hat{\mathbf{q}}_h + \hat{\mathbf{u}}_h^f \mathbf{div} \hat{\mathbf{q}}_h \right) \mathbf{d}\mathbf{x} = 0 \quad \forall \hat{\mathbf{q}}_h \in \hat{\mathbf{Q}}_{0h}, \quad (149)$$

$$\int_{\Omega} (\mathbf{div} \hat{\mathbf{p}}_h^f + \tilde{\mathbf{f}}) \hat{\mathbf{v}}_h \mathbf{d}\mathbf{x} = 0 \quad \forall \hat{\mathbf{v}}_h \in \hat{\mathbf{V}}_h. \quad (150)$$

Under the assumptions made $\tilde{\mathbf{f}} \in \mathbf{P}^0(\mathbf{T})$, the above DM problem is solvable, and

$$\begin{aligned} \|\hat{\mathbf{p}}_h^f\|_*^2 + \int_{\Omega} \hat{\mathbf{u}}_h^f \mathbf{div} \hat{\mathbf{p}}_h^f \mathbf{d}\mathbf{x} &= 0, \\ \|\hat{\mathbf{p}}_h^f\|_*^2 &\leq \|\hat{\mathbf{u}}_h^f\| \|\mathbf{div} \hat{\mathbf{p}}_h^f\| = \|\hat{\mathbf{u}}_h^f\| \|\tilde{\mathbf{f}}\| \end{aligned}$$

From here, we observe that

$$\bar{c}_1 \|\widehat{\mathbf{p}}_h^f\|^2 \leq \|\widehat{\mathbf{p}}_h^f\|_*^2 \leq \|\widehat{\mathbf{u}}_h^f\| \|\widetilde{\mathbf{f}}\|. \quad (151)$$

By (146) and (147) we conclude that for $\widehat{\mathbf{u}}_h^f$ we can find $\bar{\mathbf{q}}_h$ in $\widehat{\mathbf{Q}}_h$ such that

$$\mathbf{div} \bar{\mathbf{q}}_h + \widehat{\mathbf{u}}_h^f = \mathbf{0} \quad \text{and} \quad \|\bar{\mathbf{q}}_h\| \leq \mathbf{C} \|\widehat{\mathbf{u}}_h^f\|$$

Use $\bar{\mathbf{q}}_h$ in the first identity (149). We have,

$$\int_{\Omega} \left(\mathbf{A}^{-1} \widehat{\mathbf{p}}_h^f \cdot \bar{\mathbf{q}}_h + \widehat{\mathbf{u}}_h^f \mathbf{div} \bar{\mathbf{q}}_h \right) dx = 0$$

Thus,

$$\begin{aligned} \|\widehat{\mathbf{u}}_h^f\|^2 &= \int_{\Omega} \widehat{\mathbf{u}}_h^f \mathbf{div} \bar{\mathbf{q}}_h \leq \|\widehat{\mathbf{p}}_h^f\|_* \|\bar{\mathbf{q}}_h\|_* \leq \\ &\leq \bar{c}_2 \|\widehat{\mathbf{p}}_h^f\|_* \|\bar{\mathbf{q}}_h\| \leq \bar{c}_2 \mathbf{C} \|\widehat{\mathbf{p}}_h^f\|_* \|\widehat{\mathbf{u}}_h^f\|. \end{aligned}$$

We observe that

$$\|\widehat{\mathbf{u}}_h^f\| \leq \bar{c}_2 \mathbf{C} \|\widehat{\mathbf{p}}_h^f\|_* . \quad (152)$$

Now, we use (151) and obtain

$$\|\widehat{\mathbf{p}}_h^f\|_*^2 \leq \|\widehat{\mathbf{u}}_h^f\| \|\widetilde{\mathbf{f}}\| \leq \bar{c}_2 \mathbf{C} \|\widehat{\mathbf{p}}_h^f\|_* \|\widetilde{\mathbf{f}}\| .$$

so that

$$\bar{c}_1 \|\widehat{\mathbf{p}}_h^f\| \leq \|\widehat{\mathbf{p}}_h^f\|_* \leq \bar{c}_2 \mathbf{C} \|\widetilde{\mathbf{f}}\| . \quad (153)$$

Hence,

$$\|\widehat{\mathbf{p}}_h^f\|_{\text{div}}^2 = \|\widehat{\mathbf{p}}_h^f\|^2 + \|\text{div}\widehat{\mathbf{p}}_h^f\|^2 \leq (1 + \frac{c_2^2}{c_1^2} \mathbf{C}^2) \|\widetilde{\mathbf{f}}\|^2 . \quad (154)$$

We note that the estimates (152), (153), and (154) show that the modified DM problem is **stable**, i.e. its solutions $(\widehat{\mathbf{p}}_h^f, \widehat{\mathbf{u}}_h^f)$ are bounded by the problem data uniformly with respect to \mathbf{h} .

If replace $\widetilde{\mathbf{f}}$ by \mathbf{f} , then we can derive the same stability estimate for the functions $(\widehat{\mathbf{p}}_h, \widehat{\mathbf{u}}_h)$ that present an approximate solution of the original DM problem.

Projection estimates for fluxes

Now, we return to the projection error estimates. As we have seen

$$\| \mathbf{p} - \mathbf{p}_h \|_* \leq \inf_{\mathbf{q}_h \in \mathbf{Q}_h} \| \mathbf{p} - \mathbf{q}_h \|.$$

This estimate did not satisfy us because the set \mathbf{Q}_h is difficult to construct.

To avoid this drawback, we apply the following procedure.

Let $\boldsymbol{\eta}_h = \widehat{\mathbf{p}}_h^f + \widehat{\mathbf{q}}_h$, where $\widehat{\mathbf{q}}_h$ is *an arbitrary element of $\widehat{\mathbf{Q}}_h$* .

We have,

$$\begin{aligned} \operatorname{div} \boldsymbol{\eta}_h &= \operatorname{div} \widehat{\mathbf{p}}_h^f + \operatorname{div} \widehat{\mathbf{q}}_h = -\widetilde{\mathbf{f}} + \operatorname{div} \widehat{\mathbf{q}}_h = \\ &= \operatorname{div}(\mathbf{p} - \widehat{\mathbf{q}}_h) + \operatorname{div} \widehat{\mathbf{q}}_h = \operatorname{div} \mathbf{p} = -\mathbf{f}. \end{aligned}$$

Therefore, $\boldsymbol{\eta}_h \in \mathbf{Q}_f$

Now, we recall the projection inequality and substitute in it $\boldsymbol{\eta}_h$:

$$\| \mathbf{p} - \mathbf{p}_h \|_* \leq \| \mathbf{p} - \boldsymbol{\eta}_h \|_* = \| \mathbf{p} - \widehat{\mathbf{p}}_h^f - \widehat{\mathbf{q}}_h \|_* \leq \quad (155)$$

$$\leq \| \mathbf{p} - \widehat{\mathbf{q}}_h \|_* + \| \widehat{\mathbf{p}}_h^f \|_* . \quad (156)$$

Note that in the case considered $\mathbf{div}(\mathbf{p} - \mathbf{p}_h) = \mathbf{0}$, so that

$$\| \mathbf{p} - \mathbf{p}_h \|_{\mathbf{div}} = \| \mathbf{p} - \mathbf{p}_h \| \leq \frac{1}{\bar{c}_1} \| \mathbf{p} - \mathbf{p}_h \|_* .$$

Therefore,

$$\| \mathbf{p} - \mathbf{p}_h \|_{\mathbf{div}} \leq \frac{1}{\bar{c}_1} (\| \mathbf{p} - \widehat{\mathbf{q}}_h \|_* + \| \widehat{\mathbf{p}}_h^f \|_*) \leq$$

$$(153) \leq \frac{1}{\bar{c}_1} (\| \mathbf{p} - \widehat{\mathbf{q}}_h \|_* + \bar{c}_2 \mathbf{C} \| \tilde{\mathbf{f}} \|) .$$

But $\tilde{\mathbf{f}} = \mathbf{div}(\mathbf{p} - \widehat{\mathbf{q}}_h)$.

Thus, we arrive at the estimate

$$\begin{aligned} \|\mathbf{p} - \mathbf{p}_h\|_{\text{div}} &\leq \\ &\leq \frac{1}{\bar{c}_1} (\|\mathbf{p} - \hat{\mathbf{q}}_h\|_* + \bar{c}_2 \mathbf{C} \|\mathbf{div}(\mathbf{p} - \hat{\mathbf{q}}_h)\|) \quad \forall \hat{\mathbf{q}}_h \in \hat{\mathbf{Q}}_h. \end{aligned}$$

and, therefore,

$$\|\mathbf{p} - \mathbf{p}_h\|_{\text{div}} \leq \bar{C}_p \inf_{\hat{\mathbf{q}}_h \in \hat{\mathbf{Q}}_h} \{ \|\mathbf{p} - \hat{\mathbf{q}}_h\|_* + \|\mathbf{div}(\mathbf{p} - \hat{\mathbf{q}}_h)\| \}. \quad (157)$$

where \bar{C}_p depends on \mathbf{C} , \bar{c}_1 , and \bar{c}_2 and does not depend on \mathbf{h} .

Projection type error estimates for $\hat{u} - \hat{u}_h$

We have

$$\int_{\Omega} (\mathbf{A}^{-1} \hat{\mathbf{p}}_h \cdot \hat{\mathbf{q}}_h + \hat{u}_h \operatorname{div} \hat{\mathbf{q}}_h) \, dx = 0 \quad \forall \hat{\mathbf{q}}_h \in \hat{\mathbf{Q}}_h.$$

Since $\hat{\mathbf{Q}}_h \subset \mathbf{Q}$, we also have

$$\int_{\Omega} (\mathbf{A}^{-1} \mathbf{p} \cdot \hat{\mathbf{q}}_h + u \operatorname{div} \hat{\mathbf{q}}_h) \, dx = 0.$$

From here, we observe that

$$\int_{\Omega} (\mathbf{A}^{-1} (\hat{\mathbf{p}}_h - \mathbf{p}) \cdot \hat{\mathbf{q}}_h + (\hat{u}_h - u) \operatorname{div} \hat{\mathbf{q}}_h) \, dx = 0 \quad \forall \hat{\mathbf{q}}_h \in \hat{\mathbf{Q}}_h.$$

Denote

$$[\mathbf{u}]_{\mathbf{T}} = \frac{1}{|\mathbf{T}|} \int_{\mathbf{T}} \mathbf{u} \, dx, \quad [\mathbf{u}]_{\mathbf{h}}(\mathbf{x}) = [\mathbf{u}]_{\mathbf{T}_i} \quad \text{if } \mathbf{x} \in \mathbf{T}_i.$$

Since $\mathbf{div} \hat{\mathbf{q}}_{\mathbf{h}}$ is constant on each \mathbf{T}_i , we rewrite the relation as follows:

$$\int_{\Omega} \left(\mathbf{A}^{-1}(\hat{\mathbf{p}}_{\mathbf{h}} - \mathbf{p}) \cdot \hat{\mathbf{q}}_{\mathbf{h}} + (\hat{\mathbf{u}}_{\mathbf{h}} - [\mathbf{u}]_{\mathbf{h}}) \mathbf{div} \hat{\mathbf{q}}_{\mathbf{h}} \right) dx = 0 \quad \forall \hat{\mathbf{q}}_{\mathbf{h}} \in \hat{\mathbf{Q}}_{\mathbf{h}}.$$

Note that $[\mathbf{u}]_{\mathbf{h}} \in \hat{\mathbf{V}}_{\mathbf{h}}$ and therefore $\bar{\mathbf{u}}_{\mathbf{h}} := \hat{\mathbf{u}}_{\mathbf{h}} - [\mathbf{u}]_{\mathbf{h}} \in \hat{\mathbf{V}}_{\mathbf{h}}$. Now, we exploit the compatibility and stability conditions (146) and (147) again. For $\bar{\mathbf{u}}_{\mathbf{h}}$ one can find $\mathbf{q}'_{\mathbf{h}} \in \hat{\mathbf{Q}}_{\mathbf{h}}$ such that

$$\mathbf{div} \mathbf{q}'_{\mathbf{h}} + \bar{\mathbf{u}}_{\mathbf{h}} = \mathbf{0} \quad \text{and} \quad \|\mathbf{q}'_{\mathbf{h}}\| \leq \mathbf{C} \|\bar{\mathbf{u}}_{\mathbf{h}}\|.$$

Let us use this function \mathbf{q}'_h in the integral relation. We have

$$\int_{\Omega} \left(\mathbf{A}^{-1}(\hat{\mathbf{p}}_h - \mathbf{p}) \cdot \mathbf{q}'_h + \bar{\mathbf{u}}_h \operatorname{div} \mathbf{q}'_h \right) \mathbf{d}\mathbf{x} = 0.$$

From here, we conclude that

$$\begin{aligned} \|\bar{\mathbf{u}}_h\|^2 &= \left| \int_{\Omega} \mathbf{A}^{-1}(\hat{\mathbf{p}}_h - \mathbf{p}) \cdot \mathbf{q}'_h \right| \leq \\ &\leq \|\hat{\mathbf{p}}_h - \mathbf{p}\|_* \|\mathbf{q}'_h\|_* \leq \mathbf{C} \bar{c}_2 \|\hat{\mathbf{p}}_h - \mathbf{p}\|_* \|\bar{\mathbf{u}}_h\|. \end{aligned}$$

Thus,

$$\|\bar{\mathbf{u}}_h\| = \|[\mathbf{u}]_h - \hat{\mathbf{u}}_h\| \leq \mathbf{C} \bar{c}_2 \|\hat{\mathbf{p}}_h - \mathbf{p}\|_* .$$

We have

$$\begin{aligned}\|\mathbf{u} - \widehat{\mathbf{u}}_h\| &\leq \|\mathbf{u} - [\mathbf{u}]_h\| + \|[\mathbf{u}]_h - \widehat{\mathbf{u}}_h\| \leq \\ &\leq \|\mathbf{u} - [\mathbf{u}]_h\| + \mathbf{C} \bar{c}_2 \|\widehat{\mathbf{p}}_h - \mathbf{p}\|_*\end{aligned}$$

Note that by the definition of $[\mathbf{u}]_h$

$$\|\mathbf{u} - [\mathbf{u}]_h\| \leq \|\mathbf{u} - \mathbf{v}_h\| \quad \forall \mathbf{v}_h \in \widehat{\mathbf{V}}_h.$$

From here, we observe that

$$\|\mathbf{u} - \widehat{\mathbf{u}}_h\| \leq \mathbf{C} \bar{c}_2 \|\widehat{\mathbf{p}}_h - \mathbf{p}\|_* + \inf_{\mathbf{v}_h \in \widehat{\mathbf{V}}_h} \|\mathbf{u} - \mathbf{v}_h\|$$

Recall (156) and observe that

$$\begin{aligned}\|\mathbf{p} - \mathbf{p}_h\|_* &\leq \|\mathbf{p} - \widehat{\mathbf{q}}_h\|_* + \|\widehat{\mathbf{p}}_h^f\|_* \leq \\ &\|\mathbf{p} - \widehat{\mathbf{q}}_h\|_* + \bar{c}_2 \mathbf{C} \|\operatorname{div}(\mathbf{p} - \widehat{\mathbf{q}}_h)\|.\end{aligned}$$

Then, we arrive at the projection type error estimate for the primal variable

$$\begin{aligned} \|\mathbf{u} - \hat{\mathbf{u}}_h\| \leq & \\ & \leq \mathbf{C}_u \inf_{\hat{\mathbf{q}}_h \in \hat{\mathbf{Q}}_h} \left\{ \|\mathbf{p} - \hat{\mathbf{q}}_h\|_* + \|\operatorname{div}(\mathbf{p} - \hat{\mathbf{q}}_h)\| + \right. \\ & \left. + \inf_{\mathbf{v}_h \in \hat{\mathbf{V}}_h} \|\mathbf{u} - \mathbf{v}_h\| \right\}, \quad (158) \end{aligned}$$

where \mathbf{C}_u depends on \mathbf{C} , $\bar{\mathbf{c}}_1$, and $\bar{\mathbf{c}}_2$ and does not depend on \mathbf{h} . Estimates (157) and (158) lead to a qualified a priori convergence estimates provided that the solution possesses proper regularity.

Now, we proceed to the derivation of **functional type a posteriori error estimates** for the Primal Mixed and Dual Mixed methods.

Our analysis follows the lines of the paper
[S. Repin and A. Smolianski, RJNAMM \(2005\)](#).

Our aim is to find the difference between (\mathbf{u}, \mathbf{p}) and (\mathbf{v}, \mathbf{q}) in the respective energy norm, e.g. in the norm of the space $\mathbf{V} \times \mathbf{U}$.

A posteriori error estimates for PMM

A posteriori estimates for the mixed formulation are based on the relation that we have already derived:

$$\| \mathbf{p} - \tilde{\mathbf{q}} \|_*^2 + \| \nabla(\mathbf{u} - \mathbf{v}) \|^2 = 2(\mathbf{J}(\mathbf{v}) - \mathbf{I}^*(\tilde{\mathbf{q}})),$$

where $\tilde{\mathbf{q}} \in Q_f$ and $\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0$.

Since the difference of the functionals in the right-hand side can be estimated by the known way, we arrive at the estimate

$$\begin{aligned} \|\mathbf{p} - \tilde{\mathbf{q}}\|_*^2 + \|\nabla(\mathbf{u} - \mathbf{v})\|^2 &\leq 2(1 + \beta)\mathbf{D}(\nabla\mathbf{v}, \mathbf{y}) \\ &+ \left(1 + \frac{1}{\beta}\right) \mathbf{C}^2 \left(\|\mathbf{div}\mathbf{y} + \mathbf{f}\|^2 + \|\mathbf{y} \cdot \mathbf{n} - \mathbf{F}\|_{\partial_2\Omega}^2\right), \quad (159) \end{aligned}$$

where $\mathbf{y} \in \widehat{\mathbf{Q}}^+$, $\tilde{\mathbf{q}} \in \mathbf{Q}_f$ and $\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0$ are arbitrary functions and β is any positive number.

Thus, for the error in the primal variable we have

$$\begin{aligned} \|\nabla(\mathbf{u} - \mathbf{u}_h)\|^2 &\leq 2(1 + \beta)\mathbf{D}(\nabla\mathbf{u}_h, \mathbf{y}) \\ &\quad + \left(1 + \frac{1}{\beta}\right) \mathbf{C}^2 \left(\|\mathbf{divy} + \mathbf{f}\|^2 + \|\mathbf{y} \cdot \mathbf{n} - \mathbf{F}\|_{\partial_2\Omega}^2\right). \end{aligned} \quad (160)$$

where \mathbf{C} is a constant in the inequality

$$\|\mathbf{w}\|^2 + \|\mathbf{w}\|_{\partial_2\Omega}^2 \leq \mathbf{C}^2 \|\nabla\mathbf{w}\|^2 \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

A posteriori estimate for the dual variable

Error estimates for the dual variable in the dual energy norm $\|\cdot\|_*$ can be obtained by the arguments similar to those used above.

Let $\mathbf{q} \in \mathbf{Y}^*$ be an approximation of \mathbf{p} . For any $\tilde{\mathbf{q}} \in \mathbf{Q}_f^*$, we obtain (from the triangle inequality and Young inequalities with $\gamma > 0$)

$$\|\mathbf{q} - \mathbf{p}\|_*^2 \leq (1 + \gamma) \|\mathbf{q} - \tilde{\mathbf{q}}\|_*^2 + \left(1 + \frac{1}{\gamma}\right) \|\tilde{\mathbf{q}} - \mathbf{p}\|_*^2.$$

Use the fact that $\|\tilde{\mathbf{q}} - \mathbf{p}\|_*^2 \leq 2(\mathbf{J}(\mathbf{v}) - \mathbf{I}^*(\tilde{\mathbf{q}}))$.

Therefore, for any $\mathbf{w} \in \mathbf{V}_0 + \mathbf{u}_0$, we have

$$\begin{aligned} \|\mathbf{q} - \mathbf{p}\|_*^2 &\leq (\mathbf{1} + \gamma) \|\mathbf{q} - \tilde{\mathbf{q}}\|_*^2 + 2 \left(\mathbf{1} + \frac{\mathbf{1}}{\gamma} \right) (\mathbf{J}(\mathbf{u}) - \mathbf{I}^*(\tilde{\mathbf{q}})) \leq \\ &\leq (\mathbf{1} + \gamma) \|\mathbf{q} - \tilde{\mathbf{q}}\|_*^2 + 2 \left(\mathbf{1} + \frac{\mathbf{1}}{\gamma} \right) (\mathbf{J}(\mathbf{w}) - \mathbf{I}^*(\tilde{\mathbf{q}})) = \end{aligned}$$

Recall that

$$\mathbf{J}(\mathbf{w}) - \mathbf{I}^*(\tilde{\mathbf{q}}) \leq (\mathbf{1} + \beta) \mathbf{D}(\nabla \mathbf{w}, \mathbf{y}) + \left(\mathbf{1} + \frac{\mathbf{1}}{\beta} \right) \frac{\mathbf{1}}{2} \|\mathbf{q} - \tilde{\mathbf{q}}\|_*^2$$

so that the right-hand side is estimated by

$$(\mathbf{1} + \gamma) \left(\mathbf{1} + \frac{\mathbf{1}}{\gamma} + \frac{\mathbf{1}}{\beta\gamma} \right) \mathbf{d}_{\mathbf{Q}_f}^2(\mathbf{q}) + 2(\mathbf{1} + \beta) \left(\mathbf{1} + \frac{\mathbf{1}}{\gamma} \right) \mathbf{D}(\nabla \mathbf{w}, \mathbf{q}).$$

If $\mathbf{q} \in \mathbf{U}$, then $\mathbf{d}_{\mathbf{Q}_f}$ is given by the negative norm and we obtain

$$\begin{aligned} \|\mathbf{q} - \mathbf{p}\|_*^2 &\leq (1 + \gamma) \left(1 + \frac{1}{\gamma} + \frac{1}{\beta\gamma}\right) \|\mathbf{f} + \operatorname{div}\mathbf{q}\| + \\ &\quad + (1 + \beta) \left(1 + \frac{1}{\gamma}\right) \mathbf{D}(\nabla\mathbf{w}, \mathbf{q}). \end{aligned} \quad (161)$$

If $\mathbf{q} \in \mathbf{Q}$ then we have (for the Dirichlét problem)

$$\begin{aligned} \|\mathbf{q} - \mathbf{p}\|_*^2 &\leq (1 + \gamma) \left(1 + \frac{1}{\gamma} + \frac{1}{\beta\gamma}\right) \mathbf{C}_\Omega \|\mathbf{f} + \operatorname{div}\mathbf{q}\| + \\ &\quad + (1 + \beta) \left(1 + \frac{1}{\gamma}\right) \mathbf{D}(\nabla\mathbf{w}, \mathbf{q}). \end{aligned} \quad (162)$$

This estimate holds for any positive parameters β, γ , and any $\mathbf{w} \in \mathbf{V}_0 + \mathbf{u}_0$. Here \mathbf{w} is a "free" function in $\mathbf{V}_0 + \mathbf{u}_0$. This "freedom" can be used to make the estimate sharper.

Application of (161) to DM approximations leads (after optimization with respect to scalar parameters) to the estimate

$$\| \mathbf{p} - \mathbf{p}_h \|_* \leq \sqrt{2\mathbf{D}}^{1/2} (\nabla \mathbf{w}, \mathbf{y}) + \| \mathbf{y} - \mathbf{p}_h \|_* + 2\mathbf{C} \left(\| \operatorname{div} \mathbf{y} + \mathbf{f} \|^2 + \| \mathbf{y} \cdot \mathbf{n} - \mathbf{F} \|^2 \right)^{1/2}. \quad (163)$$

Here \mathbf{w} is an arbitrary function from $\mathbf{V}_0 + \mathbf{u}_0$ and \mathbf{y} is an arbitrary function from $\widehat{\mathbf{Q}}^+$. If $\mathbf{y} = \mathbf{A}\nabla \mathbf{u}$ and $w = u$, then the right-hand side of (163) coincides with the left-hand side, i.e. is exact in the sense that there exist such "free variables" that the inequality holds as the equality.

A directly computable upper bound of $\| \mathbf{p} - \mathbf{p}_h \|_*$ is given by (163), if we set

$$\mathbf{v} = \mathbf{u}_h, \quad \text{and} \quad \mathbf{y} = \mathcal{G}_h \mathbf{p}_h,$$

where $\mathcal{G}_h : \mathbf{Q}_h \rightarrow \widehat{\mathbf{Q}}^+$ is a certain projection operator (some examples such operators has been already discussed in the previous lectures). We have

$$\begin{aligned} \| \mathbf{p} - \mathbf{p}_h \|_* \leq & \sqrt{2\mathbf{D}}^{1/2} (\nabla \mathbf{u}_h, \mathcal{G}_h \mathbf{p}_h) + \| \mathcal{G}_h \mathbf{p}_h - \mathbf{p}_h \|_* \\ & + 2\mathbf{C} \left(\| \operatorname{div} \mathcal{G}_h \mathbf{p}_h + \mathbf{f} \|^2 + \| \mathcal{G}_h \mathbf{p}_h \cdot \mathbf{n} - \mathbf{F} \|^2 \right)^{1/2}. \end{aligned}$$

Projection from Q_h onto \widehat{Q}^+

If \mathbf{p}_h is a piecewise-constant vector field on a simplicial mesh \mathcal{T}_h , then, Raviart–Thomas elements (e.g., \mathbf{RT}^0 –elements) can be used in order to define the mapping \mathcal{G} .

Assume that the Ω has a polygonal boundary, and the latter is exactly matched by the triangulation \mathcal{T}_h . Let \mathbf{T}_i and \mathbf{T}_j be two neighboring simplexes with the common edge \mathbf{E}_{ij} . Let \mathbf{q}_h be a piecewise constant vector-valued function that has the values \mathbf{q}_i and \mathbf{q}_j on \mathbf{T}_i and \mathbf{T}_j respectively. Let \mathbf{E}_{ij} be the common edge with the unit normal \mathbf{n}_{ij} oriented from \mathbf{T}_i to \mathbf{T}_j if $i > j$.

How to define the common value $\tilde{\mathbf{q}}_{ij} \cdot \mathbf{n}_{ij}$ on \mathbf{E}_{ij} ?

One possible option is as follows:

$$\tilde{\mathbf{q}}_{ij} \cdot \mathbf{n}_{ij} = \frac{1}{2}(\mathbf{q}_i + \mathbf{q}_j) \cdot \mathbf{n}_{ij},$$

Another option is

$$\tilde{\mathbf{q}}_{ij} \cdot \mathbf{n}_{ij} = \frac{|\mathbf{T}_i| \mathbf{q}_i + |\mathbf{T}_j| \mathbf{q}_j}{|\mathbf{T}_i| + |\mathbf{T}_j|} \cdot \mathbf{n}_{ij},$$

where $|\mathbf{T}_i|$ and $|\mathbf{T}_j|$ are the areas of \mathbf{T}_i and \mathbf{T}_j . We repeat this procedure for all internal edges of \mathcal{T}_h .

If $\mathbf{E}_{i0} \in \partial_1 \Omega$, then we set $\tilde{\mathbf{q}}_{i0} \cdot \mathbf{n}_{i0} = \mathbf{q}_{i0} \cdot \mathbf{n}_{i0}$. If $\mathbf{E}_{i0} \in \partial_2 \Omega$, then

$$\tilde{\mathbf{q}}_{i0} \cdot \mathbf{n}_{i0} = \frac{1}{|\mathbf{E}_{i0}|} \int_{\mathbf{E}_{i0}} \mathbf{F} \, ds.$$

Here $|\mathbf{E}_{i0}|$ is the length of the edge \mathbf{E}_{i0} .

Thus, all the normal components $\tilde{\mathbf{q}}_{ij} \cdot \mathbf{n}_{ij}$ on internal and external edges are defined. By prolongation inside all \mathbf{T}_i , with the help of \mathbf{RT}_0 -approximations we obtain the function a piecewise affine function, which has continuous normal components at all the edges and piecewise constant normal components on $\partial\Omega$.

Therefore, we, in fact, have constructed a mapping $\mathbf{q}_h \rightarrow \tilde{\mathbf{q}}_h$ such that

$$\tilde{\mathbf{q}}_h = \mathcal{G}_h \mathbf{q}_h \in \hat{\mathbf{Q}}^+ .$$

A posteriori estimates for DMM

An a posteriori estimate for the flux $\widehat{\mathbf{p}}_h$ readily follows from the general estimate

$$\begin{aligned} \frac{1}{2} \|\mathbf{y} - \mathbf{p}\|_*^2 &\leq (1 + \gamma) \left(1 + \frac{1}{\gamma} + \frac{1}{\beta\gamma} \right) \|\mathbf{f} + \mathbf{div} \mathbf{I}\|^2 + \\ &\quad + (1 + \beta) \left(1 + \frac{1}{\gamma} \right) \mathbf{D}(\nabla \mathbf{w}, \mathbf{y}). \end{aligned}$$

We set $\mathbf{y} = \widehat{\mathbf{p}}_h \in \widehat{\mathbf{Q}}^+$. Since $\widehat{\mathbf{p}}_h$ is a piecewise polynomial function, it has a summable trace on $\partial_2 \Omega$. Then, we estimate $\|\mathbf{f} + \mathbf{div} \mathbf{I}\|$ from above in the same way as before.

Minimization with respect to γ and β leads to the estimate

$$\begin{aligned} \|\mathbf{p} - \widehat{\mathbf{p}}_h\|_* \leq & \sqrt{2}\mathbf{D}^{1/2}(\nabla \mathbf{w}, \widehat{\mathbf{p}}_h) + \\ & + 2\mathbf{C} \left(\|\operatorname{div} \widehat{\mathbf{p}}_h + \mathbf{f}\|^2 + \|\widehat{\mathbf{p}}_h \cdot \mathbf{n} - \mathbf{F}\|_{\partial_2 \Omega}^2 \right)^{1/2}, \end{aligned} \quad (164)$$

where \mathbf{w} is an arbitrary function from $\mathbf{V}_0 + \mathbf{u}_0$.

Assume that Ω is a polygonal domain decomposed into a regular collection of simplexes. If $\hat{\mathbf{p}}_h$ is constructed by means of \mathbf{RT}_0 -elements, then

$$\int_{\Omega} (\mathbf{div} \hat{\mathbf{p}}_h + \mathbf{f}) \mathbf{w}_h \, d\mathbf{x} = \mathbf{0} \quad \forall \mathbf{w}_h \in \hat{\mathbf{V}}_h \subset \hat{\mathbf{V}}, \quad (165)$$

where the subspace $\hat{\mathbf{V}}_h$ contains piecewise constant functions. Therefore, on each element \mathbf{T}_i

$$\mathbf{div} \hat{\mathbf{p}}_h = -\frac{1}{|\mathbf{T}_i|} \int_{\mathbf{T}_i} \mathbf{f} \, d\mathbf{x}. \quad (166)$$

Let us define by $[\mathbf{f}]$ the function whose values on \mathbf{T}_i coincide with the mean values of \mathbf{f} on \mathbf{T}_i . It is clear that $[\mathbf{f}] \in \widehat{\mathbf{V}}_h$. Then, we have

$$\operatorname{div} \widehat{\mathbf{p}}_h = -[\mathbf{f}] \quad \text{on every } \mathbf{T}_i.$$

Estimate (165) is valid for any approximate flux $\widehat{\mathbf{p}}_h$ from $\widehat{\mathbf{Q}}^+$. If $\widehat{\mathbf{p}}_h$ belongs to the narrower set $\widehat{\mathbf{Q}}_F$ (as, e.g., it would be in the discrete dual mixed method if $\mathbf{f} = [\mathbf{f}]$ and pure Dirichlet conditions) then the last norm in (165) would be identically zero.

It cannot, however, be expected, when $\widehat{\mathbf{p}}_h$ is constructed in the space \mathbf{RT}_0 , unless the function \mathbf{F} is a constant on $\partial_2 \Omega$.

Remark. The problem of taking into account the essential boundary condition for the flux variable

$$\hat{\mathbf{p}} \cdot \mathbf{n} = \mathbf{F} \quad \text{on} \quad \partial_2 \Omega$$

in the dual mixed method is not at all easy and, usually, leads to a non-conforming approximation $\hat{\mathbf{p}}_h$

(see, e.g., [I. Babuska and G. N. Gatica,](#)

[On the mixed finite element method with Lagrange multipliers. Numer. Meth. PDE, 2003](#)).

However, (165) still works for such (nonconforming) approximations of the flux !

One simple nonconforming version of the discrete dual method, particularly suited for the lowest-order Raviart-Thomas approximation is as follows. Instead of requiring $\hat{\mathbf{p}}_h \in \hat{\mathbf{Q}}_F$, we impose a weaker condition

$$\hat{\mathbf{p}}_h \cdot \mathbf{n}|_{\mathbf{E}_{i0}} = \frac{1}{|\mathbf{E}_{i0}|} \int_{\mathbf{E}_{i0}} \mathbf{F} \, ds \quad (167)$$

on every edge $\mathbf{E}_{i0} \in \partial_2 \Omega$.

If now we denote by $[\mathbf{F}]$ the piecewise constant function defined on the set of edges forming $\partial_2 \Omega$ and whose value on every edge $\mathbf{E}_{i0} \in \partial_2 \Omega$ is equal to the mean value of \mathbf{F} on that edge, we can write that $\hat{\mathbf{p}}_h \cdot \mathbf{n} = [\mathbf{F}]$ for all $\mathbf{E}_{i0} \in \partial_2 \Omega$.

In this case, **nonconformity errors will be automatically accounted in the functional a posteriori estimate** (165). Indeed, we obtain

$$\| \mathbf{p} - \hat{\mathbf{p}}_h \|_* \leq \sqrt{2} \mathbf{D}^{1/2}(\nabla \mathbf{v}, \hat{\mathbf{p}}_h) + 2\mathbf{C} \left(\|\mathbf{f} - [\mathbf{f}]\|^2 + \|\mathbf{F} - [\mathbf{F}]\|_{\partial_2 \Omega}^2 \right)^{1/2}. \quad (168)$$

How to choose in (168) the function $\mathbf{w} \in \mathbf{V}_0 + \mathbf{u}_0$.

The simplest way is to use the function $\hat{\mathbf{u}}_h \in \hat{\mathbf{V}}_h$ available from the solution of the discrete dual mixed problem and to construct a suitable projection operator $\mathcal{P}_h : \hat{\mathbf{V}}_h \rightarrow \mathbf{V}_0 + \mathbf{u}_0$. Again, the projection can be easily accomplished with a simple averaging.

Projection from $\hat{\mathbf{V}}_h$ onto $\mathbf{V}_0 + \mathbf{u}_0$.

In order to find $\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0$, it is sufficient to find $\mathbf{w} \in \mathbf{V}_0$ in the representation $\mathbf{v} = \mathbf{w} + \mathbf{u}_0$ (the function \mathbf{u}_0 is given). Using the computed piecewise-constant function $\hat{\mathbf{u}}_h$, we define $\mathbf{w}_h \in \mathbf{V}_0$ as follows.

We set

$$\mathbf{w}_h(\mathbf{x}_k) = \frac{\sum_{s=1}^{N_k} |\mathbf{T}_s^{(k)}| \cdot \hat{\mathbf{u}}_h|_{\mathbf{T}_s^{(k)}}}{\sum_{s=1}^{N_k} |\mathbf{T}_s^{(k)}|} - \mathbf{u}_0(\mathbf{x}_k) \quad (169)$$

for any internal node \mathbf{x}_k and when $\mathbf{x}_k \in \partial_2 \Omega$. Here $\mathbf{T}_s^{(k)}$, $s = \overline{1, N_k}$, are the elements containing the vertex \mathbf{x}_k , and we have assumed that the function \mathbf{u}_0 has a sufficient regularity, so that its point values are defined.

If the node $\mathbf{x}_k \in \partial_1 \Omega$, we simply set $\mathbf{w}_h(\mathbf{x}_k) = \mathbf{0}$.

Thus, using the nodal values of \mathbf{w}_h and the piecewise-linear continuous finite element approximation on the mesh \mathcal{T}_h we define the function

$$\mathbf{w}_h + \mathbf{u}_0 = \mathcal{P}_h \hat{\mathbf{u}}_h \in \mathbf{V}_0 + \mathbf{u}_0.$$

Hence, from (168) one obtains

$$\bar{c}_1 \|\mathbf{p} - \hat{\mathbf{p}}_h\| \leq \| \mathbf{p} - \hat{\mathbf{p}}_h \|_* \leq \sqrt{2\mathbf{D}}^{1/2} (\nabla(\mathcal{P}_h \hat{\mathbf{u}}_h), \hat{\mathbf{p}}_h) + 2\mathbf{C} \left(\|\mathbf{f} - [\mathbf{f}]\|^2 + \|\mathbf{F} - [\mathbf{F}]\|_{\partial_2 \Omega}^2 \right)^{1/2}, \quad (170)$$

which, together with the obvious relation

$$\|\mathbf{div}(\hat{\mathbf{p}} - \hat{\mathbf{p}}_h)\| = \| -\mathbf{f} - \mathbf{div} \hat{\mathbf{p}}_h \| = \|\mathbf{f} - [\mathbf{f}]\|$$

leads to the upper bound for $\|\hat{\mathbf{p}} - \hat{\mathbf{p}}_h\|_{\mathbf{div}}$:

Theorem

Let $(\hat{\mathbf{u}}, \hat{\mathbf{p}}) \in \hat{\mathbf{V}} \times \hat{\mathbf{Q}}_F$ be the exact solution of the dual mixed problem and $(\hat{\mathbf{u}}_h, \hat{\mathbf{p}}_h) \in \hat{\mathbf{V}}_h \times \hat{\mathbf{Q}}_{Fh}$ the solution of the discrete dual mixed problem with $\hat{\mathbf{Q}}_{Fh}$ being the Raviart-Thomas space \mathbf{RT}^0 .

Then, the following estimate holds true:

$$\|\hat{\mathbf{p}} - \hat{\mathbf{p}}_h\|_{\text{div}} \leq \|\mathbf{A}\nabla(\mathcal{P}_h\hat{\mathbf{u}}_h) - \hat{\mathbf{p}}_h\|_* + (2\mathbf{C} + 1)\|\mathbf{f} - [\mathbf{f}]\| + 2\mathbf{C}\|\mathbf{F} - [\mathbf{F}]\|_{\partial_2\Omega}, \quad (171)$$

where $\mathcal{P}_h : \hat{\mathbf{V}}_h \rightarrow \mathbf{V}_0 + \mathbf{u}_0$ is the projection (averaging) operator introduced above and $[\mathbf{f}]$ and $[\mathbf{F}]$ are the averaged functions.

Remark. The first and the second terms in (171), being computed elementwise, can serve as local error indicators.

A sharper estimate can be obtained by the minimization of the Majorant with respect to \mathbf{v} . Here, we can restrict ourselves to certain subspace V_h , i.e.,

$$\|\widehat{\mathbf{p}} - \widehat{\mathbf{p}}_h\|_{\text{div}} \leq \inf_{\mathbf{v}_h \in V_h} \|\mathbf{A}\nabla(\mathbf{v}_h) - \widehat{\mathbf{p}}_h\|_* + (2\mathbf{C} + 1)\|\mathbf{f} - [\mathbf{f}]\| + 2\mathbf{C}\|\mathbf{F} - [\mathbf{F}]\|_{\partial_2\Omega}. \quad (172)$$

By (160) we can also the squared norm of the error of the averaged solution $\mathcal{P}_h \hat{\mathbf{u}}_h$ using the computed flux approximation $\hat{\mathbf{p}}_h$:

$$\begin{aligned} \|\nabla(\mathbf{u} - \mathcal{P}_h \hat{\mathbf{u}}_h)\|^2 &\leq 2(\mathbf{1} + \beta) \mathbf{D}(\nabla(\mathcal{P}_h \hat{\mathbf{u}}_h), \hat{\mathbf{p}}_h) \\ &\quad + \left(\mathbf{1} + \frac{\mathbf{1}}{\beta}\right) \mathbf{C}^2(\|\mathbf{f} - [\mathbf{f}]\|^2 + \|\mathbf{F} - [\mathbf{F}]\|_{\partial_2 \Omega}^2), \end{aligned} \quad (173)$$

where $\beta > \mathbf{0}$ is an arbitrary number that can be used to minimize the right-hand side of (173) and to obtain the estimate for the norm of the error.

A sharper estimate may be obtained, if one spends some time on the minimization of the right-hand side of (173) with respect to the dual variable \mathbf{y} over some finite-dimensional subspace of $\hat{\mathbf{Q}}^+$.

Remark.

If one has the solutions of both the primal and the dual mixed problems, the flux approximation $\hat{\mathbf{p}}_h$ can be substituted into (160) to immediately yield the error estimate for the primal variable (which is the most important in the primal mixed method), while the approximation \mathbf{u}_h can be used in (171) to bring the error estimate for the dual variable (which is the most important in the dual mixed method).

Lecture 7

MIXED FEM ON DISTORTED MESHES

The Plan

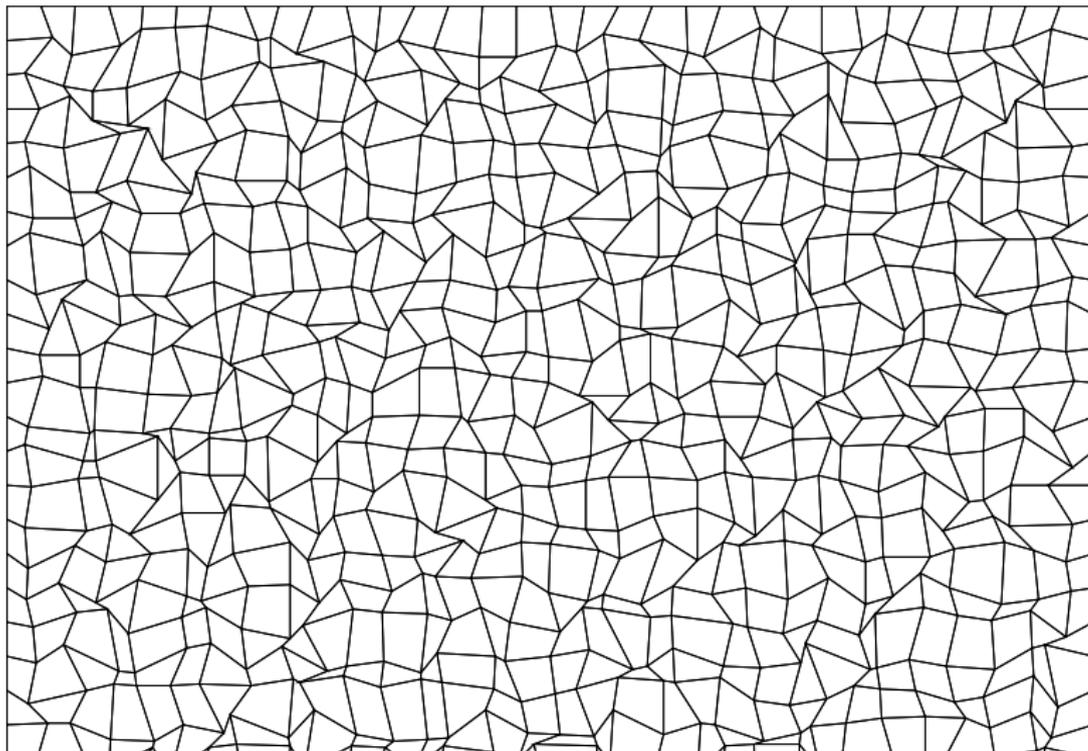
- 1 Approximations on distorted meshes
- 2 Inf-sup condition
- 3 A priori rate convergence estimates
- 4 A posteriori estimates

Distorted meshes

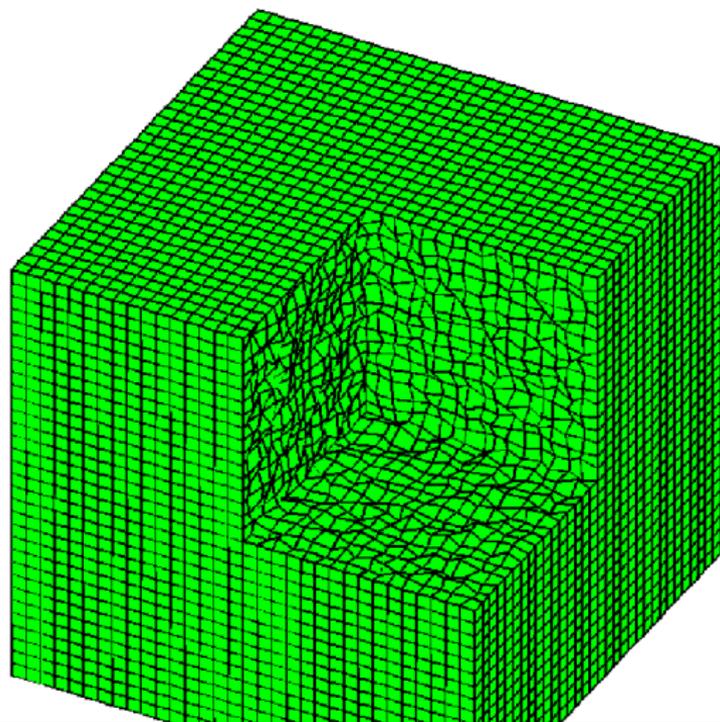
Distorted meshes may arise by several ways:

- Forcibly, due to the form of the given data;
- In the process of mesh refinement;
- In the process of adaptation to physical structure (e.g. layers of different materials)

Distorted mesh in 2D



Distorted mesh in 3D



How to work with such meshes?

Straightforward way:

Construct a regular mesh composed of standard elements that contains the given distorted one and solve the problem:

Drawback:

The problem obtained may have a very large dimension

We consider another *modus operandi* that is based upon a certain **aggregation procedure and reduction of some degrees of freedom.**

Basic problem

Principal ideas of the approach we discuss on the paradigm of the classical problem

$$-\Delta \mathbf{u} = \mathbf{f}, \quad \text{in } \Omega \quad (174)$$

$$\frac{\partial \mathbf{u}}{\partial \nu} = \mathbf{0} \quad \text{on } \partial\Omega, \quad (175)$$

Ω is a bounded connected polygonal (polyhedral) domain in \mathbb{R}^d ($d = 2, 3$), where ν is the outward unit normal vector to the boundary $\partial\Omega$ and it is assumed that

$$\mathbf{f} \in \mathbf{L}_2(\Omega), \quad \int_{\Omega} \mathbf{f} \, d\mathbf{x} = \mathbf{0}. \quad (176)$$

However, the approach can be extended to problems with the elliptic operator $\mathbf{div} \mathbf{A} \nabla$ and other types of boundary conditions.

To introduce a minimax formulation of the problem (174)–(175), we define the space of square integrable functions with zero mean

$$\mathbf{V}(\Omega) := \left\{ \mathbf{v} \in \mathbf{L}_2(\Omega) \mid \int_{\Omega} \mathbf{v} \, d\mathbf{x} = \mathbf{0} \right\},$$

for the primal variable u and the set

$$\mathbf{Q}(\Omega) := \{ \mathbf{q} \in \mathbf{H}(\Omega, \text{div}) \mid \mathbf{q} \cdot \nu = \mathbf{0} \text{ on } \partial\Omega \},$$

for the dual variable \mathbf{p} , where the boundary condition is understood in a generalized sense and $\mathbf{H}(\Omega, \text{div})$ is the Hilbert space of vector-valued functions with the scalar product and the norm defined by the relations

$$(\mathbf{p}, \mathbf{q}) := \int_{\Omega} (\mathbf{p} \cdot \mathbf{q} + \text{div} \mathbf{p} \, \text{div} \mathbf{q}) \, d\mathbf{x}, \quad \|\mathbf{p}\|_{\text{div}, \Omega} := \sqrt{(\mathbf{p}, \mathbf{p})}.$$

On $\mathbf{V} \times \mathbf{Q}$, we define the Lagrangian

$$\mathbf{L}(\mathbf{v}, \mathbf{q}) := - \int_{\Omega} \left(\mathbf{v} \operatorname{div} \mathbf{q} + \frac{1}{2} |\mathbf{q}|^2 + \mathbf{f} \mathbf{v} \right) \mathbf{d}\mathbf{x}.$$

As we have seen, the saddle-point problem: find $(\mathbf{u}, \mathbf{p}) \in \mathbf{V}(\Omega) \times \mathbf{Q}(\Omega)$ such that

$$\mathbf{L}(\mathbf{u}, \mathbf{q}) \leq \mathbf{L}(\mathbf{u}, \mathbf{p}) \leq \mathbf{L}(\mathbf{v}, \mathbf{p}) \quad \forall \mathbf{v} \in \mathbf{V}, \mathbf{q} \in \mathbf{Q}$$

has a unique solution, which satisfies the equations

$$\int_{\Omega} (\operatorname{div} \mathbf{p} + \mathbf{f}) \mathbf{w} \mathbf{d}\mathbf{x} = \mathbf{0} \quad \forall \mathbf{w} \in \mathbf{V}(\Omega), \quad (177)$$

$$\int_{\Omega} (\mathbf{p} \cdot \mathbf{q} + \mathbf{u} \operatorname{div} \mathbf{q}) \mathbf{d}\mathbf{x} = \mathbf{0} \quad \forall \mathbf{q} \in \mathbf{Q}(\Omega). \quad (178)$$

Discrete problems

Let \mathfrak{T}_h be a partitioning of Ω into polygonal (polyhedral) cells \mathbf{n}_s , $s = 1, 2, \dots, N$, $h \rightarrow 0$ as $N \rightarrow +\infty$.

We assume that

(a) partitioning is conforming in the sense that interfaces \mathbf{E}_{st} between cells \mathbf{n}_s and \mathbf{n}_t are straight segments for $d = 2$ and simply connected polygons for $d = 3$

(b) partitioning is quasiuniform and the cells are regularly shaped. The first assumption means that there exist two positive numbers α_1 and α_2 such that

$$\alpha_1 \mathbf{N}^{-d} \leq \text{diam}(\mathbf{n}_s) \leq \alpha_2 \mathbf{N}^{-d}, \quad s = 1, 2, \dots, N. \quad (179)$$

The second means that every cell \mathbf{n}_s of \mathfrak{T}_h can be partitioned into n_s regularly shaped triangles ($d = 2$), or tetrahedrons ($d = 3$), and

$$\max_s n_s \leq \text{const.}$$

We analyze samplings where the degrees of freedom represent normal components of the vector field on the edges (interfaces) between cells/subdomains.

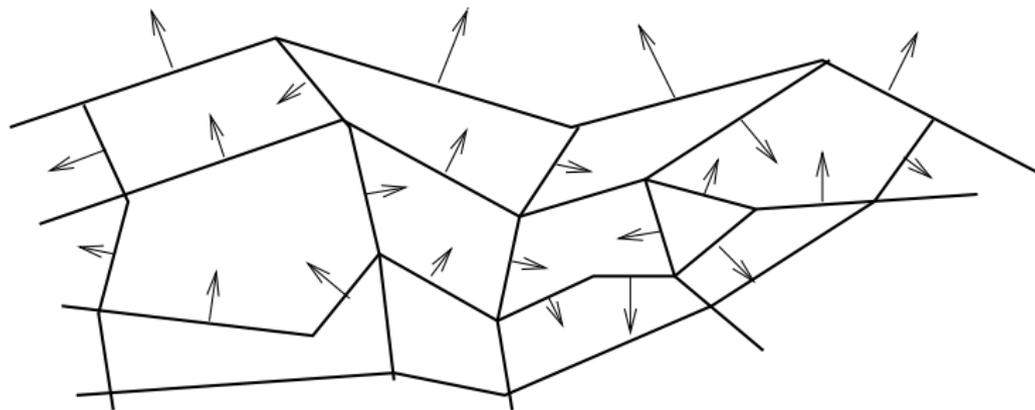
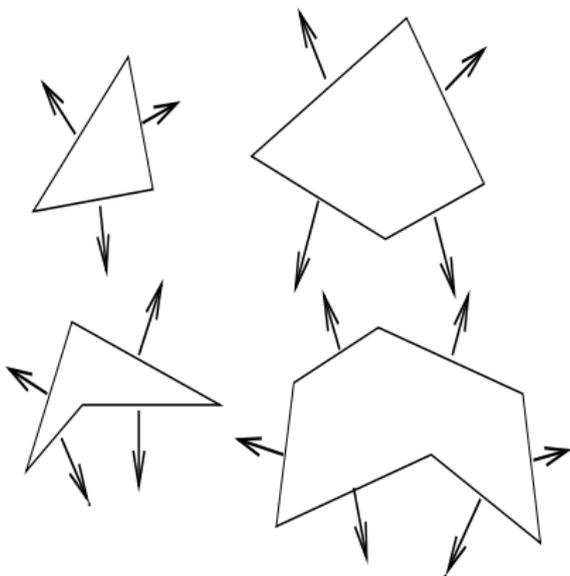
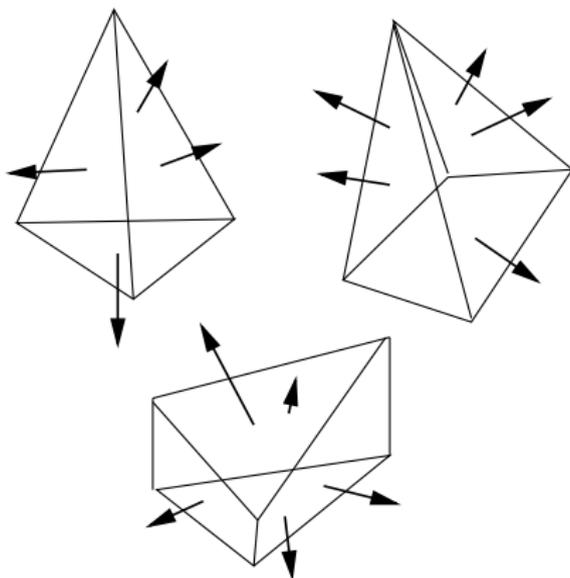


Figure: Distorted polygonal mesh

Elements of a distorted mesh in 2D



Elements of a distorted mesh in 2D



Denote by \mathbf{E}_{st} the common edge of subdomains Ω_s and Ω_t and by ν_{st} the unit normal vector to this edge oriented such that it is external to Ω_s if $s < t$ and external to Ω_t in the opposite case. The set of all edges we denote by \mathcal{E}_h . The field is approximated by the quantities \mathbf{q}_{st} defined on each $\mathbf{E}_{st} \in \mathcal{E}_h$ (see Figure 1). The value of \mathbf{q}_{st} represents the normal component of the vector field on the edge \mathbf{E}_{st} . All these quantities form the set

$$\mathbf{q}(\mathcal{T}_h) := \{ \mathbf{q}_{st} \mid \mathbf{q}_{st} \text{ is defined on } \mathbf{E}_{st} \in \mathcal{E}_h \}.$$

Besides, for each Ω_s we define a number \mathbf{V}_s . These numbers form a piecewise constant approximation of the scalar-valued function \mathbf{v}_h , i.e.

$$\mathbf{v}_h \in \mathbf{V}_h := \left\{ \mathbf{v} \in \mathbf{V}(\Omega) \mid \mathbf{v} \in \mathbf{P}^0(\Omega_s) \text{ for any } \Omega_s \right\},$$

where P^0 denotes the set of zero order polynomials.

In general, there are various finite dimensional formulations that can be established on the basis of this set of parameters. We consider the one based on the "conformity concept" and certain extension operators that transform $\mathbf{q} \in \mathbf{q}(\mathcal{T}_h)$ to a function in $H(\Omega, \text{div})$. For this purpose, we construct a linear continuous mapping $\mathbb{P}_\Omega : \mathbf{q} \rightarrow \mathbb{P}_\Omega \mathbf{q}$ such that

$$\mathcal{Q} = \mathbb{P}_\Omega \mathbf{q} \in \mathbf{H}(\Omega, \text{div}) \quad \forall \mathbf{q} \in \mathbf{q}(\mathcal{T}_h).$$

Then, the respective finite dimensional scheme readily follows from the functional formulation, if the space $\mathcal{Q}(\Omega)$ is replaced by the image space of the operator \mathbb{P}_Ω , which we denote by $\mathcal{Q}_h(\Omega)$.

To define \mathbb{P}_Ω we need first to present a suitable way of extension for a cell \mathfrak{n} .
Let

$$\mathbf{q}(\mathfrak{n}) := \{ \mathbf{q}_{\text{st}} \mid \mathbf{q}_{\text{st}} \text{ is defined on } \mathbf{E}_{\text{st}} \in \partial \mathfrak{n} \}$$

be the set of normal components of the flux on the edges of \mathfrak{n} .

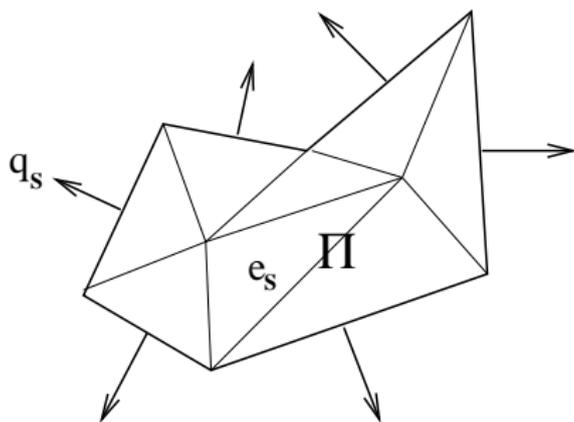


Figure: Polygonal cell with normal components of the flux

On \mathfrak{n} we define the *cell extension operator*

$$\mathbb{P}_{\mathfrak{n}} : \mathbf{q}(\mathfrak{n}) \rightarrow \mathbf{H}(\mathfrak{n}, \text{div}),$$

which maps the set of normal components of the flux given on $\partial\mathfrak{n}$ to a function $\mathcal{Q}_{\mathfrak{n}} = \mathbb{P}_{\mathfrak{n}}\mathbf{q}(\mathfrak{n})$ defined on \mathfrak{n} . This operator must be linear and satisfy the following conditions:

$$\begin{aligned} \text{(a)} \quad & \text{div } \mathcal{Q}_{\mathfrak{n}} \in \mathbf{P}^0(\mathfrak{n}), \\ \text{(b)} \quad & \mathcal{Q}_{\mathfrak{n}} \text{ is piecewise affine on } \mathfrak{n}, \\ \text{(c)} \quad & \int_{\mathfrak{n}} \text{div } \mathcal{Q}_{\mathfrak{n}} \mathbf{d}\mathbf{x} = \sum_{\mathbf{E}_{\text{st}} \in \partial\mathfrak{n}} \mathbf{q}_{\text{st}} |\mathbf{E}_{\text{st}}| \end{aligned} \tag{180}$$

Comment: why special Q_h and V_h should be used?

A cell may have a complicated form, but it is "small" so that on each sell we can approximate the pressure \mathbf{v} by a constant, i.e., the first space is

$$\mathbf{V}_h := \{\mathbf{v}_h \in \mathbf{V} = \mathbf{L}^2(\Omega) \mid \mathbf{v}_h \in \mathbf{P}^0(\mathbf{E}^i)\},$$

How then $q_h \in H(\Omega, \text{div})$ should be presented on a distorted cell ?

We use the method suggested in [Yu. Kuznetsov, S. Repin, JNM, \(2003\)](#).

It is based on the condition

$$\mathbf{div}q_h = \text{const} \quad \text{on } \mathbf{E}^i$$

Motivation

1. Approximations \mathbf{q}_h and \mathbf{v}_h are compatible in the sense that $\mathbf{div} \mathbf{q}_h \in \mathbf{V}_h$. This fact is very important in the stability and convergence analysis.
2. Also, we have approximation arguments to justify such a choice. Indeed, since $\mathbf{v}_h = \mathbf{const}$ on \mathbf{E}^i , the equation

$$\int_{\Omega} (\mathbf{div} \mathbf{p}_h + \mathbf{f}) \mathbf{v}_h \, d\mathbf{x} = 0 \quad \forall \mathbf{v}_h \in \mathbf{V}_h.$$

is equivalent to

$$\int_{\Omega} (\mathbf{div} \mathbf{p}_h + [\mathbf{f}]_{\mathbf{E}}) \mathbf{v}_h \, d\mathbf{x} = 0 \quad \forall \mathbf{v}_h \in \mathbf{V}_h.$$

where $[\mathbf{f}]_{\mathbf{E}}$ is the function taking mean values of \mathbf{f} on each cell.

Thus, in this lowest order approximations, the variability of f inside a cell is ignored and in reality we use the relation

$$\operatorname{divp} + [f]_E = 0.$$

Therefore, for the flux approximation it is natural to impose the condition:

$$\operatorname{divp}_h = \text{const} \quad \text{on } E^i$$

Certainly, \mathbb{P}_n must additionally satisfy the natural condition on that the norm $\|\mathcal{Q}_n\|_n$ must be equivalent to the norm of the vector $\mathbf{q}(n)$. This leads to certain requirements on the structure of the mesh, which are satisfied under the assumptions we have imposed on \mathcal{T}_h .

Having \mathbb{P}_{n_s} for each cell n_s , we define the *global extension operator*

$$\mathbb{P}_\Omega : \mathbf{q}(\mathcal{T}_h) \rightarrow \mathbf{Q}_h$$

by setting

$$\mathbb{P}_\Omega \mathbf{q}(\mathbf{x}) = \mathbb{P}_{n_s} \mathbf{q}(n_s)(\mathbf{x}), \quad \text{if } \mathbf{x} \in n_s$$

for any $\mathbf{q} \in \mathbf{q}(\mathcal{T}_h)$.

Consider the respective saddle point problem for the Lagrangian

$$\mathbf{L}(\mathbf{v}_h, \mathbf{q}_h) = - \int_{\Omega} \left(\mathbf{v}_h \operatorname{div} \mathbf{q}_h + \frac{1}{2} |\mathbf{q}_h|^2 + \mathbf{f} \mathbf{v}_h \right) \mathbf{d}x$$

on $\mathbf{V}_h \times \mathbf{Q}_h$. Components of the saddle point $(\mathbf{u}_h, \mathbf{p}_h)$ are solutions of two variational problems associated with the Lagrangian L . The first problem \mathcal{P}_h is to find $u_h \in V_h(\Omega)$ such that

$$\inf_{\mathbf{v}_h \in \mathbf{V}_h(\Omega)} \mathbf{J}_h(\mathbf{v}_h) = \mathbf{J}_h(\mathbf{u}_h) := \inf \mathcal{P}_h,$$

where

$$\mathbf{J}_h(\mathbf{v}_h) = \sup_{\mathbf{q}_h \in \mathbf{Q}_h} \mathbf{L}(\mathbf{v}_h, \mathbf{q}_h).$$

Another (dual) problem \mathcal{P}_h^* is to find p_h that maximizes the functional

$$I(\mathbf{q}_h) = -\frac{1}{2} \int_{\Omega} |\mathbf{q}_h|^2 \, d\mathbf{x}$$

on the set

$$Q_h(\Omega) := \left\{ \mathbf{q}_h \in \mathbf{Q}_h \mid \int_{\Omega} (\operatorname{div} \mathbf{q}_h + f) \mathbf{v}_h \, d\mathbf{x} = 0 \, \forall \mathbf{v}_h \in \mathbf{V}_h \right\}.$$

Well-posedness of the discrete problems

Under the above made assumptions on the external data and \mathfrak{T}_h both these problems \mathcal{P}_h and \mathcal{P}_h^* are stable and well-posed.

Assumption. For any $\mathbf{v}_h \in \mathbf{V}_h(\Omega)$ one can find a function $\eta_h \in Q_h(\Omega)$ such that

$$\operatorname{div} \eta_h|_{\Gamma_s} = \mathbf{v}_h|_{\Gamma_s}, \quad \mathbf{s} = 1, 2, \dots, N \quad (181)$$

$$\|\eta_h\|_{\mathbf{H}(\Omega, \operatorname{div})} \leq \mathbf{C} \|\mathbf{v}_h\|_{\Omega}, \quad (182)$$

where the constant \mathbf{C} depends on Ω .

Proposition 1. Functional J_h is coercive on V_h .

Proof

Let $\mathbf{v}_h \in \mathbf{V}_h$. For any \mathbf{v}_h , we find η_h in accordance with conditions of Proposition 1. For any positive α

$$\begin{aligned} \mathbf{J}_h(\mathbf{v}_h) &= \sup_{\mathbf{q}_h \in \mathbf{Q}_h} \mathbf{L}(\mathbf{v}_h, \mathbf{q}_h) \geq \\ &\geq \mathbf{L}(\mathbf{v}_h, -\alpha\eta_h) = \int_{\Omega} \left(\alpha \mathbf{v}_h \operatorname{div} \eta_h - \frac{\alpha^2}{2} |\eta_h|^2 - \mathbf{f} \mathbf{v}_h \right) \mathbf{d}x. \end{aligned}$$

By (181) and (182), we obtain

$$\begin{aligned} \mathbf{J}_h(\mathbf{v}_h) &\geq \sum_{s=1}^N \alpha \|\mathbf{v}_h\|_{\Pi_s}^2 - \frac{\alpha^2 \mathbf{C}^2}{2} \|\mathbf{v}_h\|_{\Omega}^2 - \int_{\Omega} \mathbf{f} \mathbf{v}_h \mathbf{d}x = \\ &= \alpha \left(\mathbf{1} - \frac{\alpha \mathbf{C}^2}{2} \right) \|\mathbf{v}_h\|_{\Omega}^2 - \int_{\Omega} \mathbf{f} \mathbf{v}_h \mathbf{d}x. \end{aligned}$$

Take $\alpha = \frac{1}{\mathbf{C}^2}$. Then

$$\begin{aligned} \mathbf{J}_h(\mathbf{v}_h) &\geq \frac{1}{2\mathbf{C}^2} \|\mathbf{v}_h\|_{\Omega}^2 - \int_{\Omega} \mathbf{f} \mathbf{v}_h \, \mathbf{d}\mathbf{x} \geq \\ &\geq \frac{1}{4\mathbf{C}^2} \|\mathbf{v}_h\|_{\Omega}^2 - \mathbf{C}^2 \|\mathbf{f}\|_{\Omega}^2 \end{aligned} \quad (183)$$

and the coercivity of J_h on V_h follows.

Proposition 2. J_h is convex and lower semicontinuous on \mathbf{V}_h .

Corollary. Since the functional J is convex, lower semicontinuous and coercive on \mathbf{V}_h , known theorems of the calculus of variations guarantee that the minimizer u_h exists (see, e.g., [?]) .

Moreover, from (183) it follows that

$$\frac{1}{4\mathbf{C}^2} \|\mathbf{u}_h\|_{\Omega}^2 \leq \mathbf{C}^2 \|\mathbf{f}\|_{\Omega}^2 + \mathbf{J}_h(\mathbf{v}_h) \quad \forall \mathbf{v}_h \in \mathbf{V}_h. \quad (184)$$

For any partition \mathfrak{T}_h the zero function belongs to \mathbf{V}_h and $\mathbf{J}_h(\mathbf{0}) = \mathbf{0}$.

Therefore, the second term in the right-hand side of (184) vanishes and we see that the norm of \mathbf{u}_h is uniformly bounded with respect to h . This means that the problems \mathcal{P}_h are stable.

It is also not difficult to justify the well-posedness of the dual problem. Indeed, Problem \mathcal{P}_h^* is a problem of maximization a quadratic functional on the affine set Q_{fh} , which has a unique solution provided that this set is not empty.

Since

$$\int_{\Omega} (\operatorname{div} \mathbf{p}_h + \mathbf{f}) \mathbf{w}_h \, \mathbf{d}\mathbf{x} = \mathbf{0} \quad \forall \mathbf{w}_h \in \mathbf{V}_h(\Omega), \quad (185)$$

$$\int_{\Omega} (\mathbf{p}_h \cdot \mathbf{q}_h + \mathbf{u}_h \operatorname{div} \mathbf{q}_h) \, \mathbf{d}\mathbf{x} = 0 \quad \forall \mathbf{q}_h \in \mathbf{Q}_h(\Omega), \quad (186)$$

we arrive at the conclusion that

$$\|\mathbf{p}_h\|_{\Omega}^2 = \int_{\Omega} \mathbf{f} \mathbf{u}_h \, \mathbf{d}\mathbf{x} \leq \|\mathbf{f}\|_{\Omega} \|\mathbf{u}_h\|_{\Omega}.$$

In view of (184), we obtain

$$\|\mathbf{p}_h\|_{\Omega}^2 \leq 2\mathbf{C}^2 \|\mathbf{f}\|_{\Omega}^2. \quad (187)$$

It is worth noting that elements of $\mathbf{Q}_{fh}(\Omega)$ satisfy the condition

$$-(\operatorname{div} \mathbf{q}_h)_{\Pi_s} = \frac{1}{|\Pi_s|} \int_{\Pi_s} \mathbf{f} dx.$$

Really, is set $w_h = 0$ on all the cells except Π_i and $w_h = c$ on that single cell. then

$$\int_{\Omega} (\operatorname{div} \mathbf{q}_h + \mathbf{f}) dx = 0 \Rightarrow [\operatorname{div} \mathbf{q}_h + \mathbf{f}]_{\Pi_i} = 0,$$

i.e.,

$$(\operatorname{div} \mathbf{q}_h)_{\Pi_s} |\Pi_s| + \int_{\Pi_s} \mathbf{f} dx = 0$$

Since $p_h \in Q_{fh}(\Omega)$, we see that

$$(\operatorname{div} \mathbf{p}_h)_{n_s} = -\frac{\mathbf{1}}{|n_s|} \int_{n_s} \mathbf{f} \, dx = \frac{\mathbf{1}}{|n_s|} \int_{n_s} \operatorname{div} \mathbf{p}. \quad (188)$$

Thus, $\operatorname{div} p_h$ **is the cell-averaging of $\operatorname{div} p$.**

From (187) and (188) it follows that

$$\|\mathbf{p}_h\|_{\text{div},\Omega} \leq \text{const.}$$

Moreover, by (188) and the property of averaged functions we conclude that if f is a smooth (or piecewise smooth) function then $\text{div} \mathbf{p}_h \rightarrow \text{div} \mathbf{p}$ a. e. in Ω .

If $f \in W_2^1(\Omega)$, then from the Poincaré inequality for the functions with zero mean we also find that

$$\|\text{div} \mathbf{p}_h - \text{div} \mathbf{p}\|_{\Omega} \rightarrow 0 \text{ as } h \rightarrow 0.$$

Other results on the convergence of \mathbf{u}_h to \mathbf{u} and \mathbf{p}_h to \mathbf{p} follow from the above stability estimates provided that the spaces \mathbf{V}_h and \mathbf{Q}_h are limit dense in \mathbf{V} and \mathbf{Q} , respectively: *for any $\mathbf{v} \in \mathbf{V}$ (resp. $\mathbf{q} \in \mathbf{Q}$) and any positive ε one can find h_ε such that for $h \leq h_\varepsilon$ \mathbf{V}_h contains a function \mathbf{v}_ε satisfying the relation $\|\mathbf{v} - \mathbf{v}_\varepsilon\|_\Omega \leq \varepsilon$ (resp. \mathbf{Q}_h contains \mathbf{q}_ε satisfying the relation $\|\mathbf{q} - \mathbf{q}_\varepsilon\|_{\text{div},\Omega} \leq \varepsilon$).*

Then, for any $\mathbf{w} \in \mathbf{V}$ and $\mathbf{q} \in \mathbf{Q}$ there exist sequences $\{\mathbf{w}_h\} \in \mathbf{V}_h$ and $\{\mathbf{q}_h\} \in \mathbf{Q}_h$ that **strongly converge** to \mathbf{w} and \mathbf{q} in \mathbf{V} and \mathbf{Q} , respectively.

Since the sequences $\{\mathbf{u}_h\}$ and $\{\mathbf{p}_h\}$ are uniformly bounded, they contain subsequences weakly converging to some limits $\tilde{\mathbf{u}} \in \mathbf{V}$ and $\tilde{\mathbf{q}} \in \mathbf{Q}$, respectively. Passing to the limit in (185)–(186) we obtain

$$\int_{\Omega} (\operatorname{div} \tilde{\mathbf{p}} + \mathbf{f}) \mathbf{w} \, \mathbf{d}\mathbf{x} = \mathbf{0} \quad \forall \mathbf{w} \in \mathbf{V}(\Omega), \quad (189)$$

$$\int_{\Omega} (\tilde{\mathbf{p}} \cdot \mathbf{q} + \tilde{\mathbf{u}} \operatorname{div} \mathbf{q}) \, \mathbf{d}\mathbf{x} = 0 \quad \forall \mathbf{q} \in \mathbf{Q}(\Omega), \quad (190)$$

what means that $\tilde{\mathbf{u}} = \mathbf{u}$ and $\tilde{\mathbf{p}} = \mathbf{p}$, i.e. the sequences converge to the solution of the basic problem.

Take sequences $\hat{\mathbf{u}}_h \rightarrow \mathbf{u}$ in \mathbf{V} and $\hat{\mathbf{p}}_h \rightarrow \mathbf{p}$ in \mathbf{Q} . We have

$$\mathbf{L}(\mathbf{u}_h, \mathbf{q}_h) \leq \mathbf{L}(\mathbf{u}_h, \mathbf{p}_h) \leq \mathbf{L}(\mathbf{v}_h, \mathbf{p}_h) \quad \forall \mathbf{v}_h \in \mathbf{V}_h, \mathbf{q}_h \in \mathbf{Q}_h. \quad (191)$$

Therefore,

$$\mathbf{J}_h(\mathbf{u}_h) = \sup_{\mathbf{q}_h \in \mathbf{Q}_h} \mathbf{L}(\mathbf{u}_h, \mathbf{q}_h) = \mathbf{L}(\mathbf{u}_h, \mathbf{p}_h), \quad (192)$$

$$\mathbf{I}(\mathbf{p}_h) = \inf_{\mathbf{v}_h \in \mathbf{V}_h} \mathbf{L}(\mathbf{v}_h, \mathbf{p}_h) = \mathbf{L}(\mathbf{u}_h, \mathbf{p}_h) \quad (193)$$

and

$$\begin{aligned} \lim_{h \rightarrow 0} I(p_h) &\geq \lim_{h \rightarrow 0} L(\hat{p}_h, u_h) = \\ &= \lim_{h \rightarrow 0} \int_{\Omega} \left(u_h \operatorname{div} \hat{p}_h - \frac{1}{2} |\hat{p}_h|^2 - f u_h \right) dx = L(u, p) = I(p). \end{aligned} \quad (194)$$

Since \mathbf{p}_h weakly converges to \mathbf{p} , we have

$$\lim_{h \rightarrow 0} \|\mathbf{p}_h\|_{\Omega} \geq \|\mathbf{p}\|_{\Omega},$$

what gives the relation

$$\lim_{h \rightarrow 0} \mathbf{I}(\mathbf{p}_h) \leq \mathbf{I}(\mathbf{p}). \quad (195)$$

From (194) and (195) we conclude that

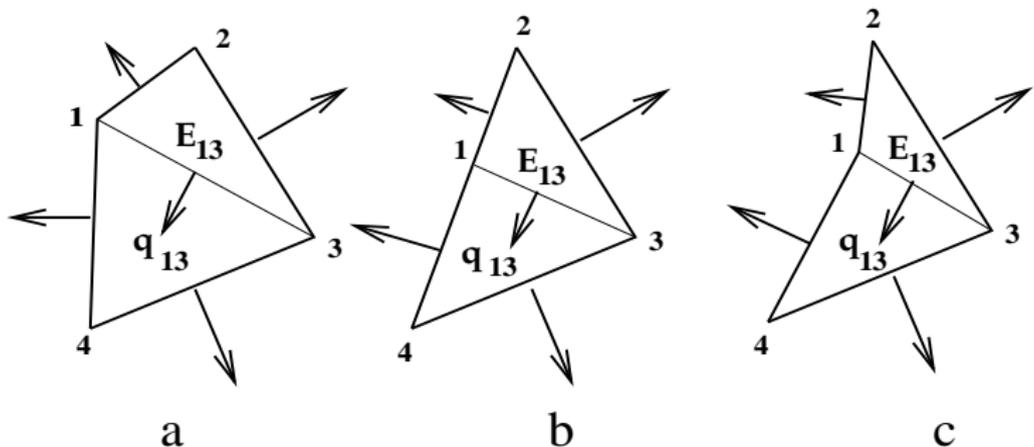
$$\lim_{h \rightarrow 0} \mathbf{I}(\mathbf{p}_h) = \mathbf{I}(\mathbf{p}). \quad (196)$$

Thus, $\|\mathbf{p}_h\|_{\Omega} \rightarrow \|\mathbf{p}\|_{\Omega}$ and, consequently, $\|\mathbf{p}_h - \mathbf{p}\|_{\Omega} \rightarrow 0$.

Examples of extension operators

Quadrilateral cells in \mathbb{R}^2 .

Quadrilateral cell π with nodes numerated 1,2,3, and 4. Edges are denoted by \mathbf{E}_{12} , \mathbf{E}_{23} , \mathbf{E}_{34} , and \mathbf{E}_{14} . Below we construct a cell extension operator \mathbb{P}_π valid for a convex, nonconvex and degenerate cells.



To construct an extended function Q_{Π} , that satisfies the conditions $Q_{\Pi}|_{\mathbf{E}_{st}} \cdot \nu_{st} = \mathbf{q}_{st}$ (where ν_{st} denotes the unit normal to \mathbf{E}_{st}) we introduce a subsidiary edge \mathbf{E}_{13} and define the respective normal component \mathbf{q}_{13} from the condition

$$\frac{1}{|\mathbf{T}_{123}|} \int_{\mathbf{T}_{123}} \operatorname{div}(Q_{\Pi}) \mathbf{d}\mathbf{x} = \frac{1}{|\mathbf{T}_{134}|} \int_{\mathbf{T}_{134}} \operatorname{div}(Q_{\Pi}) \mathbf{d}\mathbf{x}, (197)$$

which means that the divergence of the extended function Q_{Π} is constant on Π .

From (197) we obtain

$$\mathbf{q}_{13} = \rho \left(\frac{\mathbf{q}_{12}|\mathbf{E}_{12}|}{|\mathbf{T}_{123}||\mathbf{E}_{13}|} + \frac{\mathbf{q}_{23}|\mathbf{E}_{23}|}{|\mathbf{T}_{123}||\mathbf{E}_{13}|} + \frac{\mathbf{q}_{14}|\mathbf{E}_{14}|}{|\mathbf{T}_{134}||\mathbf{E}_{13}|} - \frac{\mathbf{q}_{34}|\mathbf{E}_{34}|}{|\mathbf{T}_{134}||\mathbf{E}_{13}|} \right), \quad (198)$$

where

$$\rho = \frac{|\mathbf{T}_{123}||\mathbf{T}_{134}|}{|\mathbf{T}_{123}| + |\mathbf{T}_{134}|}.$$

Now, for each of the two triangles we can construct the field using the lowest order Raviart–Thomas elements. This field meets the conditions $\operatorname{div} \mathcal{Q}_{\mathbf{n}} = \mathbf{const}$ in \mathbf{n} and

$$\int_{\mathbf{n}} \operatorname{div} \mathcal{Q}_{\mathbf{n}} \, d\mathbf{x} = \mathbf{q}_{12}|\mathbf{E}_{12}| + \mathbf{q}_{23}|\mathbf{E}_{23}| + \mathbf{q}_{34}|\mathbf{E}_{34}| + \mathbf{q}_{14}|\mathbf{E}_{14}|. \quad (199)$$

It is possible to show that the extension operator satisfies the relation

$$\|Q_{\mathbf{n}}\|_{\mathbf{n}}^2 \leq \mathbf{h}^2 \mu_2(\mathbf{n}) \sum_{E_{\text{st}} \in \partial \mathbf{n}} \mathbf{q}_{\text{st}}^2, \quad (200)$$

where μ_2 is a positive constant depending on the shape of a cell \mathbf{n} . For this purpose, at each node, we define vectors associated with certain triangle containing this node. It means that it is bounded.

Polygonal cells in \mathbb{R}^3

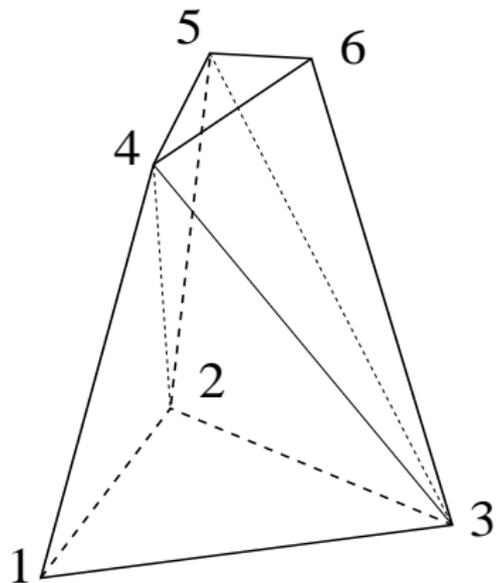


Figure: Extension inside a prism

Denote the normal components of the flux on the faces \mathbf{E}_{1245} , \mathbf{E}_{2356} , \mathbf{E}_{1346} , \mathbf{E}_{123} , and \mathbf{E}_{456} by q_1 , q_2 , q_3 , q_4 , and q_5 , respectively. Introduce two subsidiary internal faces \mathbf{E}_{234} and \mathbf{E}_{345} with normal components q' and q'' oriented inside \mathfrak{n} . Now, the prism is divided into 3 tetrahedrons T_{1234} , T_{2345} , and T_{3456} . The divergence of the extended field is defined via the Stokes theorem as

$$\operatorname{div} \mathbb{P}_{\mathfrak{n}} = \mathbf{g}_{\mathfrak{n}},$$

where

$$\mathbf{g}_{\mathfrak{n}} = \mathbf{q}_1 |\mathbf{E}_{1245}| + \mathbf{q}_2 |\mathbf{E}_{2356}| + \mathbf{q}_3 |\mathbf{E}_{1346}| + \mathbf{q}_4 |\mathbf{E}_{123}| + \mathbf{q}_5 |\mathbf{E}_{456}|.$$

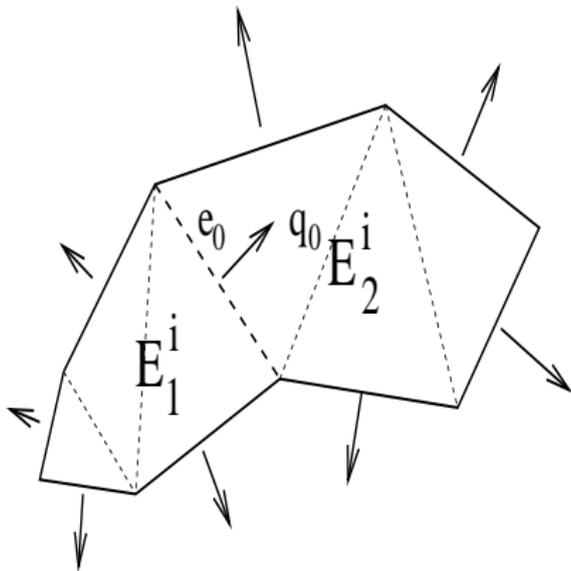
Since the condition $\text{div} \mathbb{P}_n = \mathbf{g}_n$ must hold on each of the tetrahedrons, we obtain

$$\begin{aligned} \mathbf{q}_1 |E_{124}| + \mathbf{q}' |E_{234}| + \mathbf{q}_3 |E_{134}| + \mathbf{q}_4 |E_{123}| &= \mathbf{g}_n |T_{1234}|, \\ \mathbf{q}_2 |E_{356}| + \mathbf{q}'' |E_{345}| + \mathbf{q}_3 |E_{346}| + \mathbf{q}_5 |E_{456}| &= \mathbf{g}_n |T_{3456}|. \end{aligned}$$

Thus, the numbers \mathbf{q}' and \mathbf{q}'' and the respective extension Q_n are uniquely defined.

General method of constructing \mathbb{P}_n

In more complicated cases as, e.g., for



internal fluxes are easily excluded by the $\text{div} q_h = \text{const}$ condition and the Stokes theorem.

In fact, we have constructed an **interpolation operator** π^Q on each cell. If \mathbf{q} is a vector-valued function defined on \mathbf{E}^i having summable traces, then the normal flux on $\mathbf{e}_{ij} \subset \partial\mathbf{E}^i$ can be defined as follows

$$\mathbf{q}^{ij} = \frac{1}{|\mathbf{e}_{ij}|} \int_{\mathbf{e}_{ij}} \mathbf{q} \cdot \nu_{ij} ds$$

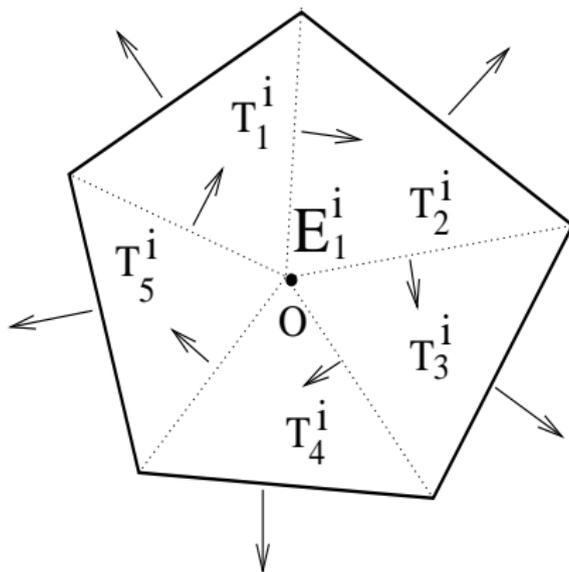
By extending the field inside \mathbf{E}^i we obtain the function \mathbf{q}_h such that

$$\begin{aligned} \mathbf{q}_h &\in \mathbf{H}(\Omega, \text{div}), \quad \text{div} \mathbf{q}_h = \text{const on } \mathbf{E}^i \\ \mathbf{q}_h &\text{ is piecewise affine in } \mathbf{E}^i \end{aligned}$$

and, thus, define the interpolation operator

$$\pi^Q : \mathcal{Q} \rightarrow \mathbf{Q}_h.$$

A cell with internal nodes



For such a sell there are many different fields that satisfy the condition

$$\mathbf{div} \mathbf{q}_h = \text{const.}$$

To avoid nonuniqueness the condition

$$\|\mathbf{q}_h\|_{E^i} \rightarrow \min$$

should be used.

We may *equivalently formulate additional condition* as follows.

Using the lowest order Raviart–Thomas finite elements we find "local" \mathcal{Q}_n satisfying the boundary condition $\mathcal{Q}_n \cdot \nu = \mathbf{q}(n) \cdot \nu$ as a solution of the discrete problem

$$\int_{\bar{n}} (\mathbf{v}_\tau \operatorname{div} \mathcal{Q}_n + \mathbf{g}_n \mathbf{v}_\tau) \mathbf{d}\mathbf{x} = \mathbf{0} \quad \forall \mathbf{v}_\tau \in \mathbf{V}_\tau, \quad (201)$$

$$\int_{\bar{n}} (\mathcal{Q}_n \cdot \mathbf{q}_\tau + \mathbf{u}_\tau \operatorname{div} \mathbf{q}_\tau) \mathbf{d}\mathbf{x} = \mathbf{0} \quad \forall \mathbf{q}_\tau \in \mathbf{Q}_{0\tau}, \quad (202)$$

where \mathbf{V}_τ is the space of piecewise constant functions and $\mathbf{Q}_{0\tau}$ is the finite element space formed by Raviart–Thomas elements contained in $H(\Omega, \operatorname{div})$ and subject to the boundary condition

$$\mathbf{q}_\tau \cdot \nu = 0 \quad \text{on } \partial n.$$

It is well known that this problem has a unique solution Q_{Ω} and the function u_{τ} is uniquely defined up to an arbitrary constant. In virtue of (201),

$$\operatorname{div} Q_{\Omega} + g_{\Omega} = 0 \quad \text{in } \Omega, \quad (203)$$

so that the requirement $\operatorname{div} Q_{\Omega} = \text{const}$ is satisfied. Since $Q_{\Omega} \in H(\Omega, \operatorname{div})$, the condition (180c) is also satisfied.

Note that the function Q_n can be also viewed as a solution of the respective discrete problem with Neumann boundary conditions, which is to maximize the functional

$$I_n(q_\tau) = -\frac{1}{2} \int_n |q_\tau|^2 dx$$

on the set of fields satisfying the condition

$$\operatorname{div} q_\tau = g_n$$

and the prescribed boundary conditions. If n is divided into a few subdomains (as, e.g., in the above example for quadrilaterals), then such a condition may uniquely define the field. However, in other cases fields with constant divergence may create a subspace. In this case, uniqueness of the extension is provided by the fact that the minimizer of the quadratic functional I_n on this subspace is uniquely defined.

Rate convergence estimates. General line

We analyze 3 different mixed approximations on distorted meshes:

A. The finest approximation on the mesh $\tilde{\mathcal{T}}_h$ when all cells are decomposed into simplexes. On simplexes the pressure field is assumed to be constant and the flux is approximated by the \mathbf{RT}^0 elements. This gives the pair of spaces $(\tilde{\mathbf{V}}_h, \tilde{\mathbf{Q}}_h)$.

B. The approximation on the mesh \mathcal{T}_h that consists of all cells. On cells the pressure field is assumed to be constant and the flux is approximated by the procedure discussed further. This gives the pair of spaces $(\mathbf{V}_h, \mathbf{Q}_h)$.

C. \mathcal{T}_h is the same, but normal fluxes are averaged on each "macroface" \mathbf{e} . This gives the pair of spaces $(\mathbf{V}_h^e, \mathbf{Q}_h^e)$.

”Referenced” problem

Consider the dual mixed formulations on $\tilde{\mathcal{T}}_h$:

Problem $\tilde{\mathcal{P}}_h$: find $(\tilde{\mathbf{p}}_h, \tilde{\mathbf{u}}_h) \in \tilde{\mathbf{Q}}_h \times \tilde{\mathbf{V}}_h$ such that

$$\int_{\Omega} (\tilde{\mathbf{p}}_h \cdot \tilde{\mathbf{q}}_h + \tilde{\mathbf{u}}_h \operatorname{div} \tilde{\mathbf{q}}_h) dx = 0 \quad \forall \tilde{\mathbf{q}}_h \in \tilde{\mathbf{Q}}_h$$

$$\int_{\Omega} (\operatorname{div} \tilde{\mathbf{p}}_h + f) \tilde{\mathbf{v}}_h dx = 0 \quad \forall \tilde{\mathbf{v}}_h \in \tilde{\mathbf{V}}_h.$$

$$\tilde{\mathbf{V}}_h := \{\tilde{\mathbf{v}}_h \in \mathbf{V} = \mathbf{L}^2(\Omega) \mid \tilde{\mathbf{v}}_h \in \mathbf{P}^0(\mathbf{T}_j^i)\},$$

$\tilde{\mathbf{Q}}_h \subset \mathbf{Q} = \mathbf{H}(\Omega, \operatorname{div})$ is constructed by \mathbf{RT}^0 elements on the mesh $\tilde{\mathcal{T}}_h$.

Properties of Problem $\tilde{\mathcal{P}}$

Problem $\tilde{\mathcal{P}}_{\mathbf{h}}$ is used as the "referenced" one. Its properties are known (see, e.g., *F. Brezzi, M. Fortin, J.E. Roberts and J.-M. Thomas*, 1991.)

1. If $\tilde{\mathcal{T}}_{\mathbf{h}}$ is regular then the inf-sup condition

$$\inf_{\tilde{\mathbf{v}}_{\mathbf{h}} \in \tilde{\mathbf{V}}_{\mathbf{h}}} \sup_{\tilde{\mathbf{q}}_{\mathbf{h}} \in \tilde{\mathbf{Q}}_{\mathbf{h}}} \frac{\int_{\Omega} \tilde{\mathbf{v}}_{\mathbf{h}} \operatorname{div} \tilde{\mathbf{q}}_{\mathbf{h}} \, dx}{\|\tilde{\mathbf{v}}_{\mathbf{h}}\| \|\tilde{\mathbf{q}}_{\mathbf{h}}\|_{\operatorname{div}, \Omega}} \geq \tilde{\mathbf{C}} \quad (204)$$

holds with a constant \mathbf{C} independent of \mathbf{h} .

2. Problem $\tilde{\mathcal{P}}_{\mathbf{h}}$ is uniquely solvable for any $\mathbf{h} > \mathbf{0}$ and

$$\tilde{\mathbf{u}}_{\mathbf{h}} \rightarrow \mathbf{u} \text{ in } \mathbf{V}(\Omega), \quad \tilde{\mathbf{p}}_{\mathbf{h}} \rightarrow \mathbf{p} \text{ in } \mathbf{Q}.$$

Moreover, for the \mathbf{RT}^0 approximations we have the standard rate convergence estimate

$$\begin{aligned} & \|\tilde{\mathbf{p}}_{\mathbf{h}} - \mathbf{p}\|_{\text{div}, \Omega} + \|\tilde{\mathbf{u}}_{\mathbf{h}} - \mathbf{u}\| \leq \\ & \leq \mathbf{C}_{\mathbf{RT}^0} \mathbf{h} (\|\mathbf{u}\|_{1, \Omega} + \|\mathbf{p}\|_{1, \Omega} + \|\mathbf{div} \mathbf{p}\|_{1, \Omega}) \end{aligned}$$

that holds provided that the exact solution is sufficiently regular.

We use the above properties
of Problem $\tilde{\mathcal{P}}_h$
in order to establish
similar properties
for
Problems \mathcal{P}_h and \mathcal{P}_h^e .

First dual mixed formulations on \mathfrak{T}_h

Problem \mathcal{P}_h : find $(\mathbf{p}_h, \mathbf{u}_h) \in \mathbf{Q}_h \times \mathbf{V}_h$ such that

$$\int_{\Omega} (\mathbf{p}_h \cdot \mathbf{q}_h + \mathbf{u}_h \operatorname{div} \mathbf{q}_h) \, dx = 0 \quad \forall \mathbf{q}_h \in \mathbf{Q}_h \quad (205)$$

$$\int_{\Omega} (\operatorname{div} \mathbf{p}_h + \mathbf{f}) \mathbf{v}_h \, dx = 0 \quad \forall \mathbf{v}_h \in \mathbf{V}_h. \quad (206)$$

$$\mathbf{V}_h := \{\mathbf{v}_h \in \mathbf{V} = \mathbf{L}^2(\Omega) \mid \mathbf{v}_h \in \mathbf{P}^0(\mathbf{E}^i), \, i = 1, 2, \dots, N\},$$

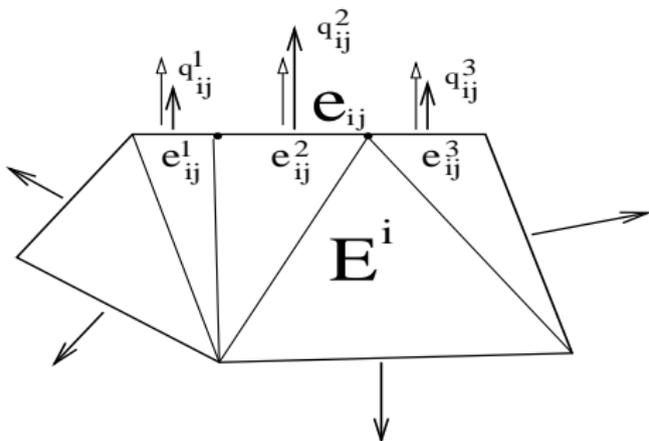
$$\mathbf{Q}_h := \{\mathbf{q}_h \in \mathbf{Q}_h \mid \operatorname{div} \mathbf{q}_h \in \mathbf{P}^0(\mathbf{E}^i), \, i = 1, 2, \dots, N\}$$

Mathematical properties of this problem are in focus of the investigation

Reduced space Q_h^e

In real life problems, analysts are often faced with meshes having "nonmatching" faces.

Certain faces of a cell may belong to one common plane and create a common boundary of one cell (macroface) denoted by Γ_h . We can reduce degrees of freedom on a macroface if replace different normal fluxes given on \mathbf{e} by one on the basis of the condition:
the value of $\int_e q_h \cdot \nu$ remains unchanged.



A macroface

By this "averaging" procedure we replace \mathbf{q}_h by $\mathbf{q}_h^e \in \mathbf{Q}_h^e$ and obtain a new space for fluxes:

$$\mathbf{Q}_h^e := \{\mathbf{q}_h \in \mathbf{Q}_h \mid \mathbf{q}_h \cdot \nu = \text{const} \quad \forall e \in \Gamma_h\}$$

Second dual mixed formulations on \mathcal{T}_h

Problem \mathcal{P}_h^e Find $(\mathbf{p}_h^e, \mathbf{u}_h^e) \in \mathbf{Q}_h^e \times \mathbf{V}_h$ such that

$$\int_{\Omega} (\mathbf{p}_h^e \cdot \mathbf{q}_h^e + \mathbf{u}_h^e \operatorname{div} \mathbf{q}_h^e) \, d\mathbf{x} = 0 \quad \forall \mathbf{q}_h^e \in \mathbf{Q}_h^e, \quad (207)$$

$$\int_{\Omega} (\operatorname{div} \mathbf{p}_h^e + \mathbf{f}) \mathbf{v}_h \, d\mathbf{x} = 0 \quad \forall \mathbf{v}_h \in \mathbf{V}_h. \quad (208)$$

It is clear that

$$\mathbf{V}_h \subset \tilde{\mathbf{V}}_h \subset \mathbf{V} \quad \text{and} \quad \mathbf{Q}_h^e \subset \mathbf{Q}_h \subset \tilde{\mathbf{Q}} \subset \mathbf{Q}.$$

Main goals

We show that if the above "aggregation–reduction" procedures are properly done and Problem $\tilde{\mathcal{P}}_h$ is stable and convergent, then

- (a) the inf–sup conditions for $(\mathbf{V}_h, \mathbf{Q}_h)$ and $(\mathbf{V}_h, \mathbf{Q}_h^e)$ are satisfied;
- (b) \mathbf{q}_h and \mathbf{q}_h^e converges to \mathbf{p} in \mathbf{Q} with the same rate as $\tilde{\mathbf{q}}_h$ and \mathbf{v}_h and \mathbf{v}_h^e converges to \mathbf{u} in \mathbf{V} with the same rate as $\tilde{\mathbf{v}}_h$;
- (c) computational errors for the approximations of all types can be explicitly controlled by the a posteriori estimates presented.

Assumptions

More precisely, we assume that

- 1** $c_1 \mathbf{h} \leq \text{diam} \mathbf{E}^i \leq c_2 \mathbf{h}$, $\forall \mathbf{E}^i \in \mathfrak{T}_h$,
- 2** Number of simplexes in a cell is from **1** to \mathbf{m}_{\max} .
- 3** $\{\tilde{\mathcal{T}}_h\}$ is regular in the usual sense:

All angles of simplexes are uniformly bounded,

Any face of a simplex is either a part of a boundary or a face of another simplex

$$\theta_1 \mathbf{h}^d \leq |\mathbf{T}_j^i| \leq \theta_2 \mathbf{h}^d, \quad \forall \mathbf{T}_j^i \in \tilde{\mathcal{T}}_h$$

If \mathbf{e} is a face of \mathbf{T}_j^i then

$$\gamma_1 \mathbf{h}^{d-1} \leq |\mathbf{e}| \leq \gamma_2 \mathbf{h}^{d-1}.$$

Logical scheme

- 1 Inf Sup in $\tilde{\mathcal{P}}_h \Rightarrow$ Inf Sup in $\mathcal{P}_h \Rightarrow$ Inf Sup in \mathcal{P}_h^e
- 2 Construct interpolation operators

$$\pi^Q : Q \rightarrow \mathbf{Q}_h \text{ and } \pi^V : \mathbf{V} \rightarrow \mathbf{V}_h$$

- 3 Establish projection type estimates
- 4 Obtain estimates for $\|\tilde{\pi}^Q \mathbf{p} - \pi^Q \mathbf{p}\|$
and $\|\tilde{\pi}^V \rho - \pi^V \rho\|$

Inf-sup condition for cell-approximations

First, we prove that for the cell-approximations Inf-Sup condition holds:

Proposition

Let $\{\tilde{\mathcal{T}}_h\}$ satisfies the above made assumptions. Then, for any $\mathbf{V}_h \in \mathbf{V}_h$ there exists a vector-valued function $\eta_h^v \in Q_h$ such that

$$-\operatorname{div} \eta_h^v = \mathbf{v}_h \quad \text{on each } \mathbf{E}^i \quad (209)$$

$$\|\eta_h^v\|_{\operatorname{div}, \Omega} \leq C \|\mathbf{v}_h\|, \quad (210)$$

where C does not depend on h and \mathbf{V}_h .

Proof.

Set $\mathbf{f} = \mathbf{v}_h$. Let $(\tilde{\mathbf{p}}_h^v, \tilde{\mathbf{u}}_h^v) \in \tilde{\mathbf{Q}}_h \times \tilde{\mathbf{V}}_h$ be the pair of functions that satisfies the system

$$\int_{\Omega} (\tilde{\mathbf{p}}_h^v \cdot \tilde{\mathbf{q}}_h + \tilde{\mathbf{u}}_h^v \operatorname{div} \tilde{\mathbf{q}}_h) \, dx = 0 \quad \forall \tilde{\mathbf{q}}_h \in \tilde{\mathbf{Q}}_h \quad (211)$$

$$\int_{\Omega} (\operatorname{div} \tilde{\mathbf{p}}_h^v + \mathbf{v}_h) \tilde{\mathbf{v}}_h \, dx = 0 \quad \forall \tilde{\mathbf{v}}_h \in \tilde{\mathbf{V}}_h. \quad (212)$$

We observe that

$$-\operatorname{div} \tilde{\mathbf{p}}_h^v = \mathbf{v}_h \quad \text{on any } \mathbf{T}_j^i.$$

Since \mathbf{v}_h is constant on \mathbf{E}^i , we arrive at the conclusion that

$$-\operatorname{div} \tilde{\mathbf{p}}_h^{\mathbf{v}} = \mathbf{v}_h \quad \text{on any } \mathbf{E}^i.$$

Now we construct the required $\eta_h^{\mathbf{v}}$ by means of $\tilde{\mathbf{p}}_h^{\mathbf{v}}$ as follows. On each cell we set normal components of $\eta_i^{\mathbf{v}}$ equal to the normal components of $\tilde{\mathbf{p}}_i^{\mathbf{v}}$. If \mathbf{E}^i contains no internal points, then (265) uniquely defines $\eta_h^{\mathbf{v}}$ by the values of normal fluxes on $\partial \mathbf{E}^i$. Thus, in this case, we simply set $\eta_h^{\mathbf{v}} = \tilde{\mathbf{p}}_h^{\mathbf{v}}$ in the cells.

If \mathbf{E}^i has an internal point, then η_h^v is defined in accordance with *minimal energy principle*, i.e.,

$$\|\eta_h^v\|_{\mathbf{E}^i}^2 = \inf_{\substack{\tau^h \in \mathfrak{N}(\mathbf{E}^i) \\ \mathbf{div} \tau_h + \mathbf{v}_h = \mathbf{0}}} \|\tau_h\|_{\mathbf{E}^i}^2. \quad (213)$$

Since $\tilde{\mathbf{p}}_h^v \in \mathfrak{N}(\mathbf{E}^i)$ and $\mathbf{div} \tilde{\mathbf{p}}_h^v + \mathbf{v}_h = \mathbf{0}$ we observe that

$$\|\eta_h^v\| \leq \|\tilde{\mathbf{p}}_h^v\|. \quad (214)$$

Set $\tilde{\mathbf{q}}_h = \tilde{\mathbf{p}}_h^v$ in (268). Then,

$$\|\tilde{\mathbf{p}}_h^v\|^2 \leq \|\tilde{\mathbf{u}}_h^v\| \|\operatorname{div} \tilde{\mathbf{p}}_h^v\| = \|\tilde{\mathbf{u}}_h^v\| \|\tilde{\mathbf{v}}_h\|. \quad (215)$$

By assumption Inf-Sup condition for $\tilde{\mathbf{Q}}_h \times \tilde{\mathbf{V}}_h$ holds! Therefore, for the function $\tilde{\mathbf{u}}_h^v$ we can find a vector-valued function $\tilde{\mathbf{q}}_h^u \in \tilde{\mathbf{Q}}_h$ such that

$$\int_{\Omega} \tilde{\mathbf{u}}_h^v \operatorname{div} \tilde{\mathbf{q}}_h^u \, dx \geq \tilde{\mathbf{C}} \|\tilde{\mathbf{u}}_h^v\| \|\tilde{\mathbf{q}}_h^u\|$$

Since

$$\int_{\Omega} (\tilde{\mathbf{p}}_h^v \cdot \tilde{\mathbf{q}}_h^u + \tilde{\mathbf{u}}_h^v \operatorname{div} \tilde{\mathbf{q}}_h^u) \, dx = 0$$

we observe that

$$\tilde{\mathbf{C}} \|\tilde{\mathbf{u}}_h^v\| \|\tilde{\mathbf{q}}_h^u\| \leq - \int_{\Omega} \tilde{\mathbf{p}}_h^v \cdot \tilde{\mathbf{q}}_h^u \, dx \leq \|\tilde{\mathbf{p}}_h^v\| \|\tilde{\mathbf{q}}_h^u\|$$

Therefore,

$$\|\tilde{\mathbf{u}}_h^{\mathbf{v}}\| \leq \frac{\mathbf{1}}{\underline{\mathbf{C}}} \|\tilde{\mathbf{p}}_h^{\mathbf{v}}\|.$$

Now, (215) leads to the estimate

$$\|\tilde{\mathbf{p}}_h^{\mathbf{v}}\| \leq \frac{\mathbf{1}}{\underline{\mathbf{C}}} \|\mathbf{v}_h\|. \quad (216)$$

Since $\eta_h^{\mathbf{v}} = \tilde{\mathbf{p}}_h^{\mathbf{v}}$, we have

$$\|\eta_h^{\mathbf{v}}\|_{\text{div}, \Omega}^2 \leq \left(\mathbf{1} + \frac{\mathbf{1}}{\underline{\mathbf{C}}^2} \right) \|\mathbf{v}_h\|^2.$$

Corollary 1. It is easy to see that under the conditions of Proposition 1, the discrete inf-sup condition for the spaces $(\mathbf{Q}_h, \mathbf{V}_h)$ holds. Indeed, take arbitrary $\mathbf{v}_h \in \mathbf{V}_h$. We have

$$\sup_{\mathbf{q}_h \in \mathbf{Q}_h} \frac{\int_{\Omega} \mathbf{v}_h \operatorname{div} \mathbf{q}_h \, \mathbf{d}\mathbf{x}}{\|\mathbf{v}_h\| \|\mathbf{q}_h\|_{\operatorname{div}, \Omega}} \geq \frac{\int_{\Omega} \mathbf{v}_h \operatorname{div} \eta_h^{\mathbf{v}} \, \mathbf{d}\mathbf{x}}{\|\mathbf{v}_h\| \|\eta_h^{\mathbf{v}}\|_{\operatorname{div}, \Omega}} \frac{\|\mathbf{v}_h\|}{\|\eta_h^{\mathbf{v}}\|_{\operatorname{div}, \Omega}} = C.$$

Thus, for $(\mathbf{V}_h, \mathbf{Q}_h)$ the inf-sup condition holds with a certain positive constant \mathbf{C} depending on $\tilde{\mathbf{C}}$.

If for cell-approximations Inf-Sup holds, then
for them PROJECTION ESTIMATE HOLDS

$$\begin{aligned} & \| \mathbf{p} - \mathbf{p}_h \|_{\text{div}, \Omega} + \| \mathbf{u} - \mathbf{u}_h \| \leq \\ & \leq \mathbf{C} \left\{ \inf_{\mathbf{q}_h \in \mathbf{Q}_h} \| \mathbf{p} - \mathbf{q}_h \|_{\text{div}, \Omega} + \inf_{\mathbf{w}_h \in \mathbf{V}_h} \| \mathbf{u} - \mathbf{w}_h \| \right\}. \end{aligned}$$

Key point in deriving a priori estimates

In projection estimates we have quantities of the type

$$\inf_{\mathbf{q}_h \in \mathbf{Q}_h} \|\mathbf{p} - \mathbf{q}_h\|, \quad \inf_{\mathbf{v}_h \in \mathbf{V}_h} \|\mathbf{u} - \mathbf{v}_h\|.$$

We estimate them by

$$\begin{aligned} \|\mathbf{p} - \pi^{\mathbf{Q}} \mathbf{p}\| &\leq \|\mathbf{p} - \tilde{\pi}^{\mathbf{Q}} \mathbf{p}\| + \|\tilde{\pi}^{\mathbf{Q}} \mathbf{p} - \pi^{\mathbf{Q}} \mathbf{p}\|, \\ \|\mathbf{u} - \pi^{\mathbf{V}} \mathbf{u}\| &\leq \|\mathbf{u} - \tilde{\pi}^{\mathbf{V}} \mathbf{u}\| + \|\tilde{\pi}^{\mathbf{V}} \mathbf{u} - \pi^{\mathbf{V}} \mathbf{u}\| \end{aligned}$$

For \mathbf{RT}^0 approximations on $\tilde{\mathcal{T}}_h$ we have the estimate

$$\|\mathbf{p} - \tilde{\pi}^{\mathbf{Q}} \mathbf{p}\| \sim \mathbf{h}, \quad \|\mathbf{u} - \tilde{\pi}^{\mathbf{V}} \mathbf{u}\| \sim \mathbf{h}.$$

Therefore, we need to estimate the difference between interpolation operators on the "fine" and "coarse" meshes.

Estimates for "pressure" interpolants.

$$\|\pi_{\mathbf{h}}^{\mathbf{V}} \mathbf{w} - \tilde{\pi}_{\mathbf{h}}^{\mathbf{V}} \mathbf{w}\| \leq \mathbf{C}_{\mathbf{V}} \mathbf{E}(\mathfrak{T}_{\mathbf{h}}, \tilde{\mathfrak{T}}_{\mathbf{h}}, \mathbf{w}), \quad (217)$$

where the constant $\mathbf{C}_{\mathbf{V}}$ does not depend on \mathbf{h} and the quantity \mathbf{E} is defined by the relations

$$\begin{aligned} \mathbf{E}(\mathfrak{T}_{\mathbf{h}}, \tilde{\mathfrak{T}}_{\mathbf{h}}, \mathbf{w}) &= \max_{i=1,2,\dots,N} \mathbf{E}_i(\mathbf{w}), \\ \mathbf{E}_i(\mathbf{w}) &= \max_{j=1,\dots,m_i} \left| [\mathbf{w}]_{\mathsf{T}_j^i} - [\mathbf{w}]_{\mathsf{E}^i} \right|. \end{aligned}$$

Indeed, take a cell \mathbf{E}^i and a simplex in it.

$$\begin{aligned}
 \int_{\mathbf{T}_j^i} (\pi_h^v \mathbf{w} - \tilde{\pi}_h^v \mathbf{w})^2 \, d\mathbf{x} &= |\mathbf{T}_j^i| \left(\frac{1}{|\mathbf{E}^i|} \int_{\mathbf{E}^i} \mathbf{w} \, d\mathbf{x} - \frac{1}{|\mathbf{T}_j^i|} \int_{\mathbf{T}_j^i} \mathbf{w} \, d\mathbf{x} \right)^2 = \\
 &= |\mathbf{T}_j^i| \left(\frac{|\mathbf{T}_j^i| \int_{\mathbf{E}^i} \mathbf{w} \, d\mathbf{x} - |\mathbf{E}^i| \int_{\mathbf{T}_j^i} \mathbf{w} \, d\mathbf{x}}{|\mathbf{E}^i| |\mathbf{T}_j^i|} \right)^2 = \\
 &= |\mathbf{T}_j^i| \left(\frac{|\mathbf{T}_j^i| \int_{\mathbf{E}^i \setminus \mathbf{T}_j^i} \mathbf{w} \, d\mathbf{x} - |\mathbf{E}^i \setminus \mathbf{T}_j^i| \int_{\mathbf{T}_j^i} \mathbf{w} \, d\mathbf{x}}{|\mathbf{E}^i| |\mathbf{T}_j^i|} \right)^2 = \\
 &= \frac{|\mathbf{T}_j^i| |\mathbf{E}^i \setminus \mathbf{T}_j^i|^2}{|\mathbf{E}^i|^2} \left([\mathbf{w}]_{\mathbf{E}^i \setminus \mathbf{T}_j^i} - [\mathbf{w}]_{\mathbf{T}_j^i} \right)^2.
 \end{aligned}$$

Now, we apply the estimate

$$\left([\mathbf{w}]_{\mathbf{E}^i \setminus \mathbf{T}_j^i} - [\mathbf{w}]_{\mathbf{T}_j^i}\right)^2 \leq 2 \left[\left([\mathbf{w}]_{\mathbf{E}^i \setminus \mathbf{T}_j^i} - [\mathbf{w}]_{\mathbf{E}^i}\right)^2 + \left([\mathbf{w}]_{\mathbf{T}_j^i} - [\mathbf{w}]_{\mathbf{E}^i}\right)^2 \right]$$

and note that

$$[\mathbf{w}]_{\mathbf{E}^i \setminus \mathbf{T}_j^i} - [\mathbf{w}]_{\mathbf{E}^i} \leq \mathbf{E}_i(\mathbf{w}).$$

By these arguments we obtain

$$\int_{\mathbf{T}_j^i} (\pi_{\mathbf{h}}^{\mathbf{V}} \mathbf{w} - \tilde{\pi}_{\mathbf{h}}^{\mathbf{V}} \mathbf{w})^2 \, \mathbf{d}\mathbf{x} \leq 2 |\mathbf{T}_j^i| \mu_i^2 \mathbf{E}_i^2(\mathbf{w}).$$

By summing over all simplexes, than over all cells we obtain (279). Since $\mu_i \leq 1$, $\mathbf{C}_{\mathbf{V}} = 2$.

Similar estimate holds true for the interpolants of fluxes

$$\|\pi_h^Q \mathbf{q} - \tilde{\pi}_h^Q \mathbf{q}\|_{\text{div}, \Omega} \leq \bar{\sigma} \mathbf{E}(\mathfrak{T}_h, \tilde{\mathfrak{T}}_h, \text{div} \mathbf{q}), \quad (218)$$

Rate convergence estimates for p_h and u_h

By the above results we arrive at *a priori rate convergence estimates* for approximations on distorted meshes (see Yu. Kuznetsov and S. Repin, JNM, 2005).

$$\begin{aligned} & \| \mathbf{p} - \mathbf{p}_h \|_{\text{div}, \Omega} + \| \mathbf{u} - \mathbf{u}_h \| \leq \\ & \leq \mathbf{C} \left[\mathbf{h} (\| \mathbf{p} \|_{1, \Omega} + \| \text{div} \mathbf{p} \|_{1, \Omega} + \| \nabla \mathbf{u} \|) + \right. \\ & \quad \left. + \mathbf{E}(\mathcal{T}_h, \tilde{\mathcal{T}}_h, \text{div} \mathbf{p}) + \mathbf{E}(\mathcal{T}_h, \tilde{\mathcal{T}}_h, \mathbf{u}) \right], \end{aligned}$$

A posteriori estimates for cell approximations

$$\|\mathbf{p} - \mathbf{p}_h\| \leq \sqrt{2} \|\nabla \mathbf{u}_h - \widehat{\mathbf{p}}_h\| + 2\mathbf{C}_\Omega \|\mathbf{div} \mathbf{p}_h + \mathbf{f}\|,$$

Since $-\mathbf{div} \mathbf{p}_h = [\mathbf{f}]_E$, we can rewrite this upper bound as

$$\|\mathbf{p} - \mathbf{p}_h\| \leq \sqrt{2} \|\nabla \mathbf{u}_h - \widehat{\mathbf{p}}_h\| + 2\mathbf{C}_\Omega \|[\mathbf{f}]_E - \mathbf{f}\|,$$

Lecture 8

A POSTERIORI ESTIMATES FOR PROBLEMS IN THE THEORY OF VISCOUS INCOMPRESSIBLE FLUIDS

Lecture plan

- **Mathematical models of viscous fluids;**
- **Stokes problem;**
- **Inf-sup condition ;**
- **A posteriori estimates for solenoidal approximations ;**
- **A posteriori estimates for non-solenoidal approximations;**
- **A posteriori estimates for problems with condition $\text{div} v = \phi$;**
- **A posteriori estimates for problems on a subspace.**
- **Bingham fluids. A posteriori estimates.**

Coordinates of particles at $t = 0$ are denoted \mathbf{X} and called **Lagrangian coordinates**. They serve as particles labels. Then, the trajectory of an individual particle is given by the relation

$$\mathbf{x} = \mathbf{x}(\mathbf{X}, t),$$

where \mathbf{x} denote the Cartesian coordinates of particles at the moment t . They are called **Euler coordinates**.

Motion equations describing evolution of a media have the form

$$\rho_0 \frac{\partial \mathbf{v}}{\partial \mathbf{t}}(\mathbf{x}, t) + \mathbf{v}_i(\mathbf{x}, t) \frac{\partial \mathbf{v}}{\partial \mathbf{x}_i}(\mathbf{x}, t) = \mathbf{div} \boldsymbol{\sigma}(\mathbf{x}, t) + \mathbf{f}(\mathbf{x}, t) \quad (219)$$

Here $\boldsymbol{\sigma}$ is the effective stress,

$\mathbf{v}_i \frac{\partial \mathbf{v}}{\partial x_i}$ is the so-called **convective term** which is often presented as

$$(\mathbf{v} \cdot \nabla) \mathbf{v}(\mathbf{x}, \mathbf{t})$$

and

$$\boldsymbol{\varepsilon}(\mathbf{v}) = \frac{1}{2} (\mathbf{v}_{i,j} + \mathbf{v}_{j,i}) \quad (220)$$

is the tensor of small strains.

In the majority of models liquids are assumed to be **incompressible**, what means that

$$\operatorname{div} \mathbf{v} = 0.$$

In view of this fact

$$\boldsymbol{\varepsilon}(\mathbf{v}) = \boldsymbol{\varepsilon}^D(\mathbf{v}).$$

Constitutive relation are usually given in the form

$$\boldsymbol{\sigma} = -\mathbf{p}\mathbb{I} + \boldsymbol{\tau}, \quad (221)$$

where \mathbf{p} is the pressure and \mathbb{I} is the unit element of $\mathbb{M}_s^{n \times n}$ and "deviatoric stress" is defined by the relation

$$\boldsymbol{\tau} \in \partial \mathbf{W}(\boldsymbol{\varepsilon}) \quad (222)$$

where $\mathbf{W} : \mathbb{M}_s^{n \times n} \mapsto \mathbb{R}_+$ is the dissipative potential and $\partial \mathbf{W}$ stands for subdifferential.

Boundary conditions

Two main types of the problems:

Flow of a fluid in a container (basin).

$$\Omega_0 = \Omega_t = \Omega, \quad t \geq 0; \quad \mathbf{v} = \mathbf{u}_0 \quad \partial\Omega \quad (223)$$

$\mathbf{v} = \mathbf{0}$ on $\partial\Omega$ is called the "adhesion" or "no-slip" condition.

Flow in an "open" domain.

In such a case, on a part of the boundary we set

$$\boldsymbol{\sigma}(\mathbf{x}, \mathbf{t}) \boldsymbol{\nu}(\mathbf{x}, \mathbf{t}) = \mathbf{F}(\mathbf{x}, \mathbf{t}), \quad \mathbf{x} \in \partial\Omega_t, \quad \mathbf{t} \geq 0 \quad (224)$$

Initial condition.

$$\mathbf{v}(\mathbf{x}, \mathbf{0}) = \boldsymbol{\varphi}(\mathbf{x}) \quad \mathbf{x} \in \Omega_0 \quad (225)$$

Now the task is to find \mathbf{v} , \mathbf{p} , $\boldsymbol{\tau}$ such that

$$\rho_0 \frac{\partial \mathbf{v}}{\partial \mathbf{t}} + \mathbf{v}_i \frac{\partial \mathbf{v}}{\partial \mathbf{x}_i} - \mathbf{div} \boldsymbol{\sigma} = \mathbf{f} \quad (226)$$

$$\mathbf{div} \mathbf{v} = 0 \quad (\mathbf{x}, \mathbf{t}) \in \mathbf{Q} := \Omega \times (0, +\infty) \quad (227)$$

$$\boldsymbol{\sigma} = -\mathbf{p}\mathbb{I} + \boldsymbol{\tau}, \quad \boldsymbol{\tau} \in \partial \mathbf{W}(\varepsilon) \quad (228)$$

$$\mathbf{v}(\mathbf{x}, \mathbf{t}) = \mathbf{0} \quad (\mathbf{x}, \mathbf{t}) \in \partial_1 \Omega \times (0, +\infty) \quad (229)$$

$$\boldsymbol{\sigma} \boldsymbol{\nu} = \mathbf{F} \quad (\mathbf{x}, \mathbf{t}) \in \partial_2 \Omega \times (0, +\infty) \quad (230)$$

$$\mathbf{v}(\mathbf{x}, 0) = \varphi(\mathbf{x}) \quad \mathbf{x} \in \Omega$$

Bingham type fluids

$$\mathbf{W}(\boldsymbol{\varepsilon}) = \frac{\mu}{\mathbf{m}} |\boldsymbol{\varepsilon}|^{\mathbf{m}} + \mathbf{k}_* |\boldsymbol{\varepsilon}| \quad (231)$$

A particular case with $\mathbf{k}_* = \mathbf{0}$ and $\mathbf{m} = 2$ leads to

$$\boldsymbol{\sigma} = -p\mathbb{I} + \frac{\partial \mathbf{W}}{\partial \boldsymbol{\varepsilon}} = -p\mathbb{I} + \mu \boldsymbol{\varepsilon},$$

Then

$$\operatorname{div} \boldsymbol{\sigma} = -\nabla p + \mu \Delta \mathbf{v}$$

and we arrive at the Navier–Stokes system

$$\begin{cases} \rho_0 \frac{\partial \mathbf{v}}{\partial t} + \mathbf{v}_i \frac{\partial \mathbf{v}}{\partial x_i} - \mu \Delta \mathbf{v} = \mathbf{f} + \nabla p \\ \operatorname{div} \mathbf{v} = 0 \end{cases} \quad (232)$$

If $k_* > 0$ then it is typical that the stagnant zones with $\boldsymbol{\varepsilon}(\mathbf{v}) \equiv \mathbf{0}$ arise.

Power law models

$$\mathbf{W}(\varepsilon) = \mu_\infty |\varepsilon|^2 + \mu_0 \left(\mathbf{1} + |\varepsilon|^2 \right)^{\frac{p}{2}} \quad (233)$$

$$\mathbf{W}(\varepsilon) = \mu_\infty |\varepsilon|^2 + \mu_0 |\varepsilon|^p \quad (234)$$

Here $\mathbf{p} \in]\mathbf{1}, +\infty[$, $\mu_\infty \geq 0$, $\mu_0 \geq 0$, $\mu_\infty \mu_0 \neq 0$.

Powel–Euring models

$$\mathbf{W}(\varepsilon) = \mu_\infty |\varepsilon|^2 + \mu_1 |\varepsilon| \ln(\mathbf{1} + |\varepsilon|) \quad (235)$$

General system

In all the cases we have the system

$$\begin{cases} \rho_0 \frac{\partial \mathbf{v}}{\partial t} + \mathbf{v}_i \frac{\partial \mathbf{v}}{\partial x_i} - \frac{\partial \mathbf{W}}{\partial \boldsymbol{\varepsilon}} (\boldsymbol{\varepsilon}(\mathbf{v})) = \mathbf{f} + \nabla \mathbf{p} \\ \mathbf{div} \mathbf{v} = \mathbf{0} \end{cases} \quad (236)$$

In case of slow motions the term $\mathbf{v} \cdot (\nabla \mathbf{v})$ is usually neglected. Then, we arrive at simplified models among which the most known is the Stokes model.

Spaces of solenoidal functions

$$\mathbf{J}^\infty(\Omega) = \{\mathbf{u} \in \mathbf{C}_0^\infty(\Omega), \operatorname{div} \mathbf{u} = \mathbf{0}, \operatorname{supp} \mathbf{u} \subset \Omega\},$$

$$\mathring{\mathbf{J}}(\Omega) := \text{closure of } \mathbf{J}^\infty \text{ in the topology of } \mathbf{L}_2(\Omega, \mathbb{R}^d),$$

$$\mathring{\mathbf{J}}_2^1(\Omega) := \text{closure of } \mathbf{J}^\infty \text{ in the topology of } \mathbf{H}^1(\Omega, \mathbb{R}^d).$$

Navier–Stokes equation

At present Navier–Stokes problem dominates among the models describing the behavior of viscous incompressible fluids. It is to find $\mathbf{u}(\mathbf{x}, \mathbf{t}) \in \mathring{\mathbf{J}}_2^1(\Omega)$ and $\mathbf{p}(\mathbf{x}, \mathbf{t}) \in \mathring{\mathbf{L}}_2(\Omega)$ such that

$$\mathbf{u}_t - \nu \Delta \mathbf{u} + \mathbf{div}(\mathbf{u} \otimes \mathbf{u}) = \mathbf{f} - \nabla \mathbf{p} \quad \text{in } \Omega, ,$$

$$\mathbf{u}(\mathbf{x}, \mathbf{0}) = \varphi(\mathbf{x}),$$

$$\mathbf{u} = \mathbf{u}_0 \quad \text{on } \Gamma_D$$

$$\varepsilon(\mathbf{u}) \cdot \boldsymbol{\nu} + \mathbf{p}\boldsymbol{\nu} = \mathbf{g}_N \quad \text{on } \Gamma_N.$$

From the mathematical point of view NS is still a mystery. Existence of a unique solution in 3D is not yet proved even for $(0, T] \times \mathbb{R}^n$.

It is known that for sufficiently regular solenoidal $\phi(\mathbf{x})$ there exists a weak **Leray-Hopf** solution, i.e., a function

$$\mathbf{u} \in \mathbf{L}^\infty(\mathbf{0}, \mathbf{T}; \mathbf{L}^2(\mathbb{R}^n)) \cap \mathbf{L}^2(\mathbf{0}, \mathbf{T}; \mathbf{H}^1(\mathbb{R}^n))$$

Proving (or presenting a contr–example) of that NS equation possesses in $(0, T] \times \mathbb{R}^n$ a smooth solution provided that initial data are sufficiently regular forms one of the **Millennium Prize Problems** stated by the Clay Mathematical Institute.

FROM THE INTRODUCTION TO THE THIRD MILLENNIUM PRIZE PROBLEM:
"...Although these (NS) equations were written down in the 19th Century, our understanding of them remains minimal. The challenge is to make substantial progress toward a mathematical theory which will unlock the secrets hidden in the Navier-Stokes equations."

However, discrete (semidiscrete) analogs of NS equation are actively used in the Mathematical Modeling.

For example:

$$\frac{\mathbf{u}^k - \mathbf{u}^{k-1}}{\Delta t} - \nu \Delta \mathbf{u}^k + \operatorname{div}(\mathbf{u}^{k-1} \otimes \mathbf{u}^{k-1}) = \mathbf{f} - \nabla \mathbf{p}^k \quad \text{in } \Omega,$$

$$\operatorname{div} \mathbf{u}^k = 0$$

and

$$\frac{\mathbf{u}^k - \mathbf{u}^{k-1}}{\Delta t} - \nu \Delta \mathbf{u}^k + \operatorname{div}(\mathbf{u}^{k-1} \otimes \mathbf{u}^k) = \mathbf{f} - \nabla \mathbf{p}^k \quad \text{in } \Omega,$$

$$\operatorname{div} \mathbf{u}^k = 0.$$

For these reasons, an intensive investigations has been devoted to *linearizations* of NS equations, in particulat to [Stokes](#)

$$\begin{aligned} \mathbf{u}_t - \nu \Delta \mathbf{u} &= \mathbf{f} - \nabla p && \text{in } \Omega, \\ \mathbf{u}(\mathbf{x}, 0) &= \varphi(\mathbf{x}), \\ \mathbf{u} &= \mathbf{u}_0 && \text{on } \Gamma_D \\ \varepsilon(\mathbf{u}) \cdot \boldsymbol{\nu} + p\boldsymbol{\nu} &= \mathbf{g}_N && \text{on } \Gamma_N \end{aligned}$$

and [Oseen](#)

$$\begin{aligned} \mathbf{u}_t - \nu \Delta \mathbf{u} + \operatorname{div}(\mathbf{a} \otimes \mathbf{u}) &= \mathbf{f} - \nabla p && \text{in } \Omega, \\ \mathbf{u}(\mathbf{x}, 0) &= \varphi(\mathbf{x}), && \operatorname{div} \mathbf{a} = \mathbf{0}, \\ \mathbf{u} &= \mathbf{u}_0 && \text{on } \Gamma_D \\ \varepsilon(\mathbf{u}) \cdot \boldsymbol{\nu} + p\boldsymbol{\nu} &= \mathbf{g}_N && \text{on } \Gamma_N \end{aligned}$$

problems.

NS equations with rotation

In certain cases (e.g., ocean modeling) NS equations should take into account Earth rotation. Then, an additional term arises and the equation comes in the form

$$\mathbf{u}_t - \nu \Delta \mathbf{u} + (\mathbf{u} \cdot \nabla) \mathbf{u} + \boldsymbol{\omega} \times \mathbf{u} = \mathbf{f} - \nabla p \quad \text{in } \Omega,$$

where

$$\boldsymbol{\omega} = |\omega| \mathbf{e}_3.$$

is the parameter of the rotation intensity.

Computer simulation methods

A significant part of the difficulties arising in the process of solving such problems is related to the **incompressibility condition $\operatorname{div} u = 0$** .

Typically, this condition is taken into account by projecting of a discrete solution to the set of solenoidal fields or by introducing appropriate penalty terms. A detailed exposition of the numerical methods can be found, e.g., in the works of (see the List of Literature) [J. Chen](#), [A. Chorin](#), [W. E](#) and [J. G. Liu](#), [M. Feistauer](#), [M. Ganzburger](#), [V. Girault](#), [G. Heywood](#), [R. Rannacher](#), [P. A. Raviart](#), [R. Temam](#).

Stationary problems are often solved by passing to a minimax formulation and using the so-called mixed approximations for the velocity and pressure fields (see, e.g., [F. Brezzi and J. Duglas](#), [F. Brezzi and M. Fortin](#)).

Linear models in the theory of fluids

Classical formulation of the Stokes problem: find a vector-valued function \mathbf{u} (velocity) and a scalar-valued function \mathbf{p} (pressure) that satisfy the relations

$$-\nu \Delta \mathbf{u} = \mathbf{f} - \nabla \mathbf{p} \quad \text{in } \Omega, \quad (237)$$

$$\operatorname{div} \mathbf{u} = 0 \quad \text{in } \Omega, \quad (238)$$

$$\mathbf{u} = \mathbf{u}_0 \quad \text{on } \partial\Omega, \quad (239)$$

where \mathbf{u}_0 is a given function such that $\operatorname{div} \mathbf{u}_0 = 0$.

Here Δ denotes the Laplacian of a vector field:

$$\Delta \mathbf{u} = (\Delta u_i) \mathbf{e}_i,$$

\mathbf{e}_i are the cartesian unit vectors.

Multiply (237) by a function $\mathbf{v} \in \mathbf{J}^\infty$. Then,

$$\int_{\Omega} \nabla \mathbf{p} \cdot \mathbf{v} \, d\mathbf{x} = \int_{\Omega} \mathbf{p} \operatorname{div} \mathbf{v} \, d\mathbf{x} = 0,$$

and we arrive at the relation

$$\nu((\mathbf{u}, \mathbf{v})) = (\mathbf{f}, \mathbf{v}) \quad \forall \mathbf{v} \in \mathbf{J}^\infty(\Omega), \quad (240)$$

where $((\mathbf{u}, \mathbf{v})) = \int_{\Omega} \boldsymbol{\varepsilon}(\mathbf{u}) : \boldsymbol{\varepsilon}(\mathbf{v}) \, d\mathbf{x}$. Above identity can be extended by

continuity to the whole $\mathbf{J}_2^1(\Omega)$.

Weak formulation of the Stokes problem

. Find $\mathbf{u} \in \mathring{\mathbf{J}}_2^1(\Omega)$, such that $\mathbf{u} = \mathbf{0}$ on $\partial\Omega$ and

$$\mu((\mathbf{u}, \mathbf{v})) = (\mathbf{f}, \mathbf{v}) \quad \forall \mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega)(\Omega), \quad (241)$$

This formulation is related to the energy functional

$$\mathbf{I}(\mathbf{v}) = \mu \int_{\Omega} |\nabla \mathbf{v}|^2 \, d\mathbf{x} - 2 \int_{\Omega} \mathbf{f} \mathbf{v} \, d\mathbf{x} \quad (242)$$

where $|\nabla \mathbf{v}|^2 = |\nabla \mathbf{v}_1|^2 + \dots + |\nabla \mathbf{v}_n|^2$.

Theorem

Generalized solution of the Stokes problem minimizes the functional \mathbf{I} over $\mathring{\mathbf{J}}_2^1(\Omega)$.

Proof. Assume that

$$\mu((\mathbf{u}, \mathbf{v})) = (\mathbf{f}, \mathbf{v}) \quad \forall \mathbf{v} \in \mathbf{V}.$$

Since $\|\mathbf{u} - \mathbf{v}\|^2 \geq 0$, we have

$$\mu \|\mathbf{u}\|^2 + \mu \|\mathbf{v}\|^2 - 2\mu((\mathbf{u}, \mathbf{v})) \geq 0. \quad (243)$$

On the other hand,

$$\mathbf{I}(\mathbf{u}) = \mu \|\mathbf{u}\|^2 - 2(\mathbf{f}, \mathbf{u}) = \mu \|\mathbf{u}\|^2 - 2\mu((\mathbf{u}, \mathbf{u})) = -\mu \|\mathbf{u}\|^2$$

Therefore, we rewrite (243) as follows:

$$\mu \|\mathbf{v}\|^2 - 2\mu((\mathbf{u}, \mathbf{v})) \geq \mathbf{I}(\mathbf{u}),$$

Therefore,

$$\mathbf{I}(\mathbf{v}) \geq \mathbf{I}(\mathbf{u}) \quad \forall \mathbf{v} \in \mathbf{V}. \quad (244)$$

Assume that \mathbf{u} is the minimizer, then for any $\mathbf{v} \in \mathbf{V}$

$$\begin{aligned} \mathbf{I}(\mathbf{u} + \lambda \mathbf{v}) &\geq \mathbf{I}(\mathbf{u}), \\ \mu \|\mathbf{u}\|^2 + 2\lambda\mu((\mathbf{u}, \mathbf{v})) + \mu\lambda^2 \|\mathbf{v}\|^2 - 2(\mathbf{f}, \mathbf{u} + \lambda \mathbf{v}) &\geq \\ &\geq \mu \|\mathbf{u}\|^2 - 2(\mathbf{f}, \mathbf{u}), \end{aligned}$$

and, consequently,

$$\mu((\mathbf{u}, \mathbf{v})) - (\mathbf{f}, \mathbf{v}) \geq -\frac{\lambda}{2} \|\mathbf{v}\|^2 \implies \mu((\mathbf{u}, \mathbf{v})) - (\mathbf{f}, \mathbf{v}) = 0 \quad \forall \mathbf{v} \in \mathbf{V}.$$

Nomenclature

Next, $\mathbf{W} := \mathbf{W}_2^1(\Omega, \mathbb{R}^d)$ and $\Sigma := \mathbf{L}_2(\Omega, \mathbb{M}^{d \times d})$, where $\mathbb{M}^{d \times d}$ is the space of symmetric $\mathbf{d} \times \mathbf{d}$ matrixes (tensors), whose scalar product is denoted by two dots.

\mathbf{W}_0 is a subspace of \mathbf{W} that contains functions with zero traces on $\partial\Omega$.

$\mathbf{W}_0 + \mathbf{u}_0$ contains functions of the form $\mathbf{w} + \mathbf{u}_0$, where $\mathbf{w} \in \mathbf{V}_0$.

Analogously, $\mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0$ contains functions of the form $\mathbf{w} + \mathbf{u}_0$, $\mathbf{w} \in \mathring{\mathbf{J}}_2^1(\Omega)$.

The operator $\varepsilon(\mathbf{v}) := \frac{1}{2}(\nabla \mathbf{v} + (\nabla \mathbf{v})^T)$ acts from \mathbf{W} to Σ .

We will also use the Hilbert space $\Sigma_{\text{div}}(\Omega)$, which is a subspace of Σ that contains tensor-valued functions τ , such that $\text{div}\tau \in \mathbf{L}_2$. The scalar product in this space is defined by the relation

$$(\tau, \eta) := \int_{\Omega} (\tau : \eta + \text{div}\tau \cdot \text{div}\eta) \, dx.$$

As before, $\mathring{\mathbf{L}}_2(\Omega)$ denotes the space of square summable functions with zero mean. Henceforth, we assume that

$$\mathbf{f} \in \mathbf{L}_2(\Omega, \mathbb{R}^d), \quad \mathbf{u}_0 \in \mathbf{W}_2^1(\Omega, \mathbb{R}^d),$$

Generalized solution can be defined by the **integral identity**. It is a function $\mathbf{u} \in \mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0$ that meets the relation

$$\int_{\Omega} \nu \nabla(\mathbf{u}) : \nabla(\mathbf{v}) \, d\mathbf{x} = \int_{\Omega} \mathbf{f} \cdot \mathbf{v} \, d\mathbf{x} \quad \forall \mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega). \quad (245)$$

It is well known that \mathbf{u} exists and unique and can be viewed as the minimizer of the functional

$$I(\mathbf{v}) = \int_{\Omega} \left(\frac{\nu}{2} |\nabla(\mathbf{v})|^2 - \mathbf{f} \cdot \mathbf{v} \right) d\mathbf{x}$$

on the set $\mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0$. Thus, the problem

$$\inf_{\mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0} I(\mathbf{v})$$

presents a **variational formulation** of the Stokes problem.

Existence of a minimizer follows from known properties of convex lower semicontinuous functionals.

In addition, the Stokes problem can be presented in a **minimax form**.

Let $\mathbf{L} : (\mathbf{W}_0 + \mathbf{u}_0) \times \mathring{\mathbf{L}}_2(\Omega) \rightarrow \mathbb{R}$ be defined as follows:

$$\mathbf{L}(\mathbf{v}, \mathbf{q}) = \int_{\Omega} \left(\frac{\nu}{2} |\nabla \mathbf{v}|^2 - \mathbf{f} \cdot \mathbf{v} - \mathbf{q} \operatorname{div} \mathbf{v} \right) dx.$$

Now, \mathbf{u} and \mathbf{p} are defined as a saddle-point that satisfies the relations

$$\mathbf{L}(\mathbf{u}, \mathbf{q}) \leq \mathbf{L}(\mathbf{u}, \mathbf{p}) \leq \mathbf{L}(\mathbf{v}, \mathbf{p}) \quad \forall \mathbf{v} \in \mathbf{W}_0 + \mathbf{u}_0, \mathbf{p} \in \mathring{\mathbf{L}}_2(\Omega).$$

Extension of solenoidal fields and related results

First, we recall some basic results that has been established when the **solvability of the Stokes** problem was investigated. Works of O.A. Ladyzhenskaya have made a considerable and widely acknowledged contribution to the mathematical theory of viscous incompressible fluids.

The first principal result states that a solenoidal field can be extended inside a domain such that the norm of the extended field is subject to the norm of the boundary trace (see O.A. Ladyzhenskaya *Mathematical problems in the dynamics of a viscous incompressible fluid*. Nauka, Moscow, 1970 and

O.A Ladyzhenskaya and V.A. Solonnikov Some problems of vector analysis, and generalized formulations of boundary value problems for the Navier-Stokes equation, *Zap. Nauchn. Sem. Leningrad. Otdel. Mat. Inst. Steklov. (LOMI)*, 59(1976), 81–116, 256).

Lemma 1.

For any vector-valued function $\mathbf{a} \in \mathbf{W}_2^{1/2}(\partial\Omega)$ satisfying the condition $\int_{\partial\Omega} \mathbf{a} \cdot \boldsymbol{\nu} \, d\mathbf{x} = \mathbf{0}$ there exists a function $\bar{\mathbf{u}} \in \mathbf{W}_0$ such that $\operatorname{div} \bar{\mathbf{u}} = 0$ and

$$\|\nabla \bar{\mathbf{u}}\| \leq \kappa_1(\Omega) \|\mathbf{a}\|_{1/2, \partial\Omega}, \quad (246)$$

where $\kappa_1(\Omega)$ is a positive constant that depends on Ω .

This lemma implies another proposition, which is of great importance for the analysis of problems defined on solenoidal fields.

Lemma 2

For any $\mathbf{f} \in \mathring{\mathbf{L}}_2(\Omega)$ there exists a function $\bar{\mathbf{u}} \in \mathbf{W}_0$ satisfying the relation $\mathbf{div} \bar{\mathbf{u}} = \mathbf{f}$ and the condition

$$\|\nabla \bar{\mathbf{u}}\| \leq \kappa_2(\Omega) \|\mathbf{f}\|, \quad (247)$$

where $\kappa_2(\Omega)$ is a positive constant that depends on Ω .

Lemma 2 implies several important corollaries that we discuss below.

Inf-Sup condition

Lemma 2 is related to the inequality known in the literature as the **Inf-Sup**– or **LBB (Ladyzhenskaya–Babuška–Brezzi)**–condition that reads: **there exists a positive constant C such that**

$$\inf_{\substack{\phi \in \mathring{L}_2(\Omega) \\ \phi \neq 0}} \sup_{\substack{\mathbf{w} \in \mathbf{W}_0 \\ \mathbf{w} \neq 0}} \frac{\int_{\Omega} \phi \operatorname{div} \mathbf{w} \, dx}{\|\phi\| \|\nabla \mathbf{w}\|} \geq C. \quad (248)$$

Inf-Sup condition (248) was established in the papers by I. Babuška The finite element method with Lagrangian multipliers, *Numer. Math.*, 20(1973) and F. Brezzi, On the existence, uniqueness and approximation of saddle-point problems arising from Lagrange multipliers, *R.A.I.R.O., Annal. Numer.*, 8 (1974). They used its discrete analogs for proving the convergence of finite-dimensional approximations in various problems related to the theory of viscous incompressible fluids.

Lemma 2 implies LBB condition

By Lemma 2, any $\phi \in \mathring{\mathbf{L}}_2(\Omega)$ has a counterpart function $\mathbf{v}_\phi \in \mathbf{W}_0$ that meets the conditions

$$\mathbf{divv}_\phi = \phi, \quad \|\nabla \mathbf{v}_\phi\| \leq \kappa_2(\Omega) \|\phi\|.$$

In this case,

$$\sup_{\mathbf{v} \in \mathbf{W}_0, \mathbf{w} \neq 0} \frac{\int_\Omega \phi \mathbf{divv} \, \mathbf{d}\mathbf{x}}{\|\nabla \mathbf{v}\| \|\phi\|} \geq \frac{\int_\Omega \phi \mathbf{divv}_\phi \, \mathbf{d}\mathbf{x}}{\|\nabla \mathbf{v}_\phi\| \|\phi\|} = \frac{\|\phi\|}{\|\nabla \mathbf{v}_\phi\|} \geq \frac{1}{\kappa_2(\Omega)}$$

and, consequently, Inf-Sup condition holds with

$$\mathbf{C} = \frac{1}{\kappa_2(\Omega)}.$$

It is easy to observe that the Inf-Sup condition can be presented in the form

$$\sup_{\substack{\mathbf{w} \in \mathbf{W}_0 \\ \mathbf{w} \neq \mathbf{0}}} \frac{\int_{\Omega} \mathbf{p} \operatorname{div} \mathbf{w} \, dx}{\|\nabla \mathbf{w}\|} \geq \mathbf{C} \|\mathbf{p}\| \quad \text{for all } \mathbf{p} \in \mathring{\mathbf{L}}_2(\Omega).$$

We may consider the expression in the left-hand side of the above inequality as the norm of $\nabla \mathbf{p}$ in the space topologically dual to \mathbf{W}_0 , namely

$$\|\nabla \mathbf{p}\| := \sup_{\mathbf{w} \in \mathbf{W}_0} \frac{\langle \nabla \mathbf{p}, \mathbf{w} \rangle}{\|\nabla \mathbf{w}\|}.$$

Then, we arrive to the Nečas inequality.

Nečas inequality

$$\|\mathbf{p}\| \leq \kappa_2 \|\nabla \mathbf{p}\| \quad \forall \mathbf{p} \in \mathring{\mathbf{L}}_2(\Omega), \quad (249)$$

A simple proof of the Nečas inequality for domains with Lipschitz boundaries can be found in the paper by

J. Bramble. *A proof of the inf-sup condition for the Stokes equations on Lipschitz domains*, Math. Models Methods Appl. Sci. **13** (2003), no. 3, 361–371.

In the later paper, it is also shown that the well-known Korn's inequality follows from Inf-Sup condition.

Constants \mathbf{C} and κ_2 play an important role in the numerical analysis of the Stokes problem as well as in the theoretical one.

Existence of a saddle point

Existence of a saddle point of $\mathbf{L}(\mathbf{v}, \mathbf{q})$ follows from Lemma 2 and known results of the minimax theory. In a simplified version these results reads:

Lagrangian $\mathbf{L}(\mathbf{v}, \mathbf{q})$ possess a saddle point provided that

- (a) it is convex and continuous with respect to the first variable and concave and continuous with respect to the second one;**
- (b) for a certain $\bar{\mathbf{q}}$ the functional $\mathbf{v} \mapsto \mathbf{L}(\mathbf{v}, \bar{\mathbf{q}})$ is coercive (or the set of admissible \mathbf{v} is compact);**
- (c) or a certain $\bar{\mathbf{v}}$ the functional $\mathbf{q} \mapsto -\mathbf{L}(\bar{\mathbf{v}}, \mathbf{q})$ is coercive (or the set of admissible \mathbf{q} is compact.)**

Since

$$\mathbf{J}(\mathbf{v}) = \sup_{\mathbf{q} \in \Sigma} \mathbf{L}(\mathbf{v}, \mathbf{q}) \geq \mathbf{L}(\bar{\mathbf{q}}, \mathbf{v}),$$

we observe that (b) means that $\mathbf{J}(\mathbf{v})$ is coercive. Analogously, (c) means that the functional $-\mathbf{I}(\mathbf{q})$, where

$$\mathbf{I}(\mathbf{q}) = \inf_{\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0} \mathbf{L}(\mathbf{v}, \mathbf{q}) \leq \mathbf{L}(\mathbf{q}, \bar{\mathbf{v}}),$$

is coercive.

In other words, for a continuous convex-concave Lagrangian existence of a saddle point mainly depends on the coercivity properties of the two dual functionals generated by it.

Let us apply these results to the Stokes problem. It is easy to see that for any $\mathbf{q} \in \mathring{\mathbf{L}}_2(\Omega)$ the mapping

$$\mathbf{v} \mapsto \mathbf{L}(\mathbf{v}, \mathbf{q}) = \int_{\Omega} \left(\frac{\nu}{2} |\nabla \mathbf{v}|^2 - \mathbf{f} \cdot \mathbf{v} - \mathbf{q} \operatorname{div} \mathbf{v} \right) \mathbf{d}\mathbf{x}.$$

is convex and continuous (in \mathbf{W}) and there exists an element $\bar{\mathbf{q}} \in \mathring{\mathbf{L}}_2(\Omega)$ (e.g., $\bar{\mathbf{q}} = 0$) such that $\mathbf{L}(\mathbf{v}, \bar{\mathbf{q}}) \rightarrow +\infty$ if $\|\mathbf{v}\|_{\mathbf{V}} \rightarrow +\infty$. The mapping $\mathbf{q} \mapsto \mathbf{L}(\mathbf{v}, \mathbf{q})$ is affine and continuous (in $\mathring{\mathbf{L}}_2(\Omega)$) for any $\mathbf{v} \in \mathbf{V}$. Therefore, existence of a saddle point is guaranteed provided that the coercivity condition

$$\lim_{\|\mathbf{q}\| \rightarrow +\infty} \inf_{\mathbf{v} \in \mathbf{W}_0 + \mathbf{u}_0} \mathbf{L}(\mathbf{v}, \mathbf{q}) = -\infty \quad (250)$$

is established. By Lemma 2 we can prove this fact.

Consider the functional

$$\mathbf{I}(\mathbf{q}) := \inf_{\mathbf{v} \in \mathbf{W}_0 + \mathbf{u}_0} \mathbf{L}(\mathbf{v}, \mathbf{q})$$

and the variational problem

$$\mathbf{I}(\mathbf{p}) = \sup_{\mathbf{q} \in \mathring{\mathbf{L}}_2(\Omega)} \mathbf{I}(\mathbf{q}) \quad (251)$$

for the **pressure** function. Note that the functional \mathbf{I} has no explicit integral-type form and is defined as a supremum–functional. The solvability of this problem follows from the coercivity condition (250). To prove (250) we apply Lemma 2.

Coercivity of the variational problem for the pressure function

Indeed, by Lemma 2 for any $\mathbf{q} \in \mathring{\mathbf{L}}_2(\Omega)$ we find $\mathbf{v}_q \in \mathbf{W}_0$ such that

$$\operatorname{div} \mathbf{v}_q = \mathbf{q} \quad \text{and} \quad \|\nabla \mathbf{v}_q\| \leq \kappa_2 \|\mathbf{q}\|.$$

Take $\mathbf{v} = \mu \mathbf{v}_q + \mathbf{u}_0$ and recall that $\operatorname{div} \mathbf{u}_0 = \mathbf{0}$. Then,

$$\begin{aligned} \inf_{\mathbf{v} \in \mathbf{W}_0 + \mathbf{u}_0} \mathbf{L}(\mathbf{v}, \mathbf{q}) &\leq \int_{\Omega} \left(\frac{\nu}{2} |\nabla(\mu \mathbf{v}_q + \mathbf{u}_0)|^2 - \mathbf{f} \cdot (\mu \mathbf{v}_q + \mathbf{u}_0) - \mathbf{q} \operatorname{div}(\mu \mathbf{v}_q + \mathbf{u}_0) \right) \mathbf{d}\mathbf{x} \leq \\ &\leq \int_{\Omega} \left(\frac{\nu}{2} |\nabla \mathbf{u}_0|^2 - \mathbf{f} \cdot \mathbf{u}_0 \right) \mathbf{d}\mathbf{x} + \mu (\nu \|\nabla \mathbf{u}_0\| + \mathbf{C}_{\Omega} \|\mathbf{f}\|) \|\nabla \mathbf{v}_q\| + \\ &\quad + \frac{\nu \mu^2}{2} \|\nabla \mathbf{v}_q\|^2 - \mu \|\mathbf{q}\|^2 \leq \int_{\Omega} \left(\frac{\nu}{2} |\nabla \mathbf{u}_0|^2 - \mathbf{f} \cdot \mathbf{u}_0 \right) \mathbf{d}\mathbf{x} + \\ &\quad + \mu (\nu \|\nabla \mathbf{u}_0\| + \mathbf{C}_{\Omega} \|\mathbf{f}\|) \kappa_2 \|\mathbf{q}\| + \mu \left(\frac{\nu \mu \kappa_2^2}{2} - 1 \right) \|\mathbf{q}\|^2, \end{aligned}$$

where \mathbf{C}_{Ω} is a constant in the Friederichs inequality.

We see that

$$\begin{aligned} \mathbf{I}(\mathbf{q}) \leq & c_1(\mathbf{u}_0, \mathbf{f}, \nu) + \mu(\nu \|\nabla \mathbf{u}_0\| + \mathbf{C}_\Omega \|\mathbf{f}\|) \kappa_2 \|\mathbf{q}\| + \\ & + \mu \left(\frac{\nu \mu \kappa_2^2}{2} - 1 \right) \|\mathbf{q}\|^2. \end{aligned}$$

Set here $\mu = \frac{1}{\nu \kappa_2^2}$. Then

$$\inf_{\mathbf{v} \in \mathbf{W}_0 + \mathbf{u}_0} \mathbf{L}(\mathbf{v}, \mathbf{q}) \leq c_1 + c_2 \|\mathbf{q}\| - \frac{1}{2\nu \kappa_2^2} \|\mathbf{q}\|^2 \rightarrow -\infty \quad \text{as } \|\mathbf{q}\| \rightarrow +\infty.$$

Thus, we observe that the constant κ_2 arises in the quadratic term that provides the required coercivity property of the pressure functional.

Estimates of the distance to the set of solenoidal fields

Now we are concerned with the estimates of the **distance between a function $\widehat{\mathbf{v}} \in \mathbf{H}^1$ and the space of solenoidal functions**.

Estimates in \mathbf{L}_2 -norm. An estimate of the distance between $\widehat{\mathbf{v}}$ and the space

$$\mathbf{J}_2^1(\Omega) := \left\{ \mathbf{v} \in \mathbf{W}_2^1(\Omega) \mid \operatorname{div} \mathbf{v} = \mathbf{0} \right\}$$

in \mathbf{L}_2 -norm follow from the solvability of the Dirichlet problem for the Laplace operator. It is as follows:

$$\inf_{\mathbf{v}_0 \in \mathbf{J}_2^1} \|\widehat{\mathbf{v}} - \mathbf{v}_0\| \leq \mathbf{C}_F \|\operatorname{div} \widehat{\mathbf{v}}\|,$$

where \mathbf{C}_F is the constant in the Friederichs inequality.

Proof. Indeed, since the problem

$$\Delta\phi = \mathbf{f},$$

has a solution $\phi \in \mathring{\mathbf{W}}_2^1(\Omega)$ for any $\mathbf{f} \in \mathbf{L}_2(\Omega)$, we conclude that for any \mathbf{f} there exists $\mathbf{v}_f = \nabla\phi$ such that

$$\mathbf{div}\mathbf{v}_f = \mathbf{f} \quad \text{and} \quad \|\mathbf{v}_f\| \leq \mathbf{C}_F\|\mathbf{f}\|.$$

Set $\mathbf{f} = \mathbf{div}\hat{\mathbf{v}}$. Then,

$$\mathbf{div}(\mathbf{v}_f - \hat{\mathbf{v}}) = \mathbf{0},$$

so that $\mathbf{v}_0 = \mathbf{v}_f - \hat{\mathbf{v}}$ belongs to \mathbf{J}_2^1 and we observe that

$$\|\hat{\mathbf{v}} - \mathbf{v}_0\| \leq \mathbf{C}_F\|\mathbf{div}\hat{\mathbf{v}}\|.$$

Estimate in H^1 -norm. Let now $\hat{\mathbf{v}} \in \mathring{\mathbf{H}}^1$. Set $\mathbf{f} = \mathbf{div}\hat{\mathbf{v}}$. Since

$$\int_{\Omega} \mathbf{div}\hat{\mathbf{v}} \, dx = \int_{\partial\Omega} \mathbf{v} \cdot \boldsymbol{\nu} \, ds = 0,$$

we see that $\mathbf{f} \in \mathring{\mathbf{L}}_2(\Omega)$. Then, by Lemma 2, one can find $\mathbf{u}_f \in \mathbf{W}_0$ such that

$$\mathbf{div}\mathbf{u}_f = \mathbf{div}\hat{\mathbf{v}}, \quad \text{and} \quad \|\nabla\mathbf{u}_f\| \leq \kappa_2(\Omega)\|\mathbf{div}\hat{\mathbf{v}}\|.$$

In other words, there exists a solenoidal field $\mathbf{w}_0 = (\hat{\mathbf{v}} - \mathbf{u}_f) \in \mathbf{W}_0$ such that

$$\|\nabla(\hat{\mathbf{v}} - \mathbf{w}_0)\| = \|\nabla\hat{\mathbf{u}}_f\| \leq \kappa_2(\Omega)\|\mathbf{div}\hat{\mathbf{v}}\|.$$

This fact can be presented in another form

$$\inf_{\mathbf{v} \in \mathring{\mathbf{J}}_2^0(\Omega)} \|\nabla(\hat{\mathbf{v}} - \mathbf{v})\| \leq \kappa_2(\Omega) \|\mathbf{div} \hat{\mathbf{v}}\|. \quad (252)$$

Thus, for the functions with zero traces the distance to $\mathring{\mathbf{J}}_2^0(\Omega)$ in a strong norm is also measured via $\|\mathbf{div} \hat{\mathbf{v}}\|$, but with a different factor: $\kappa_2(\Omega)$.

Comments on the value of \mathbf{C}

Note that \mathbf{C} can be estimated throughout the constant \mathbf{C}_F and the constant \mathbf{C}_P in the Poincare inequality. Indeed,

$$\mathbf{C} = \inf_{\mathbf{q} \in \mathring{\mathbf{L}}_2(\Omega), \mathbf{q} \neq 0} \mathcal{E}(\mathbf{q}),$$

$$\mathcal{E}(\mathbf{q}) = \sup_{\mathbf{w} \in \mathbf{W}_0, \mathbf{w} \neq 0} \frac{\int_{\Omega} \mathbf{q} \operatorname{div} \mathbf{w} \, dx}{\|\mathbf{q}\| \|\nabla \mathbf{w}\|}.$$

For $\mathbf{q} \in \tilde{\mathbf{W}}(\Omega) := \mathring{\mathbf{L}}_2(\Omega) \cap \mathbf{W}_2^1(\Omega)$ we have

$$\begin{aligned} \mathcal{E}(\mathbf{q}) &= \sup_{\mathbf{w} \in \mathbf{W}_0, \mathbf{w} \neq 0} \frac{\int_{\Omega} \nabla \mathbf{q} \cdot \mathbf{w} \, dx}{\|\mathbf{q}\| \|\nabla \mathbf{w}\|} \leq \frac{\|\nabla \mathbf{q}\|}{\|\mathbf{q}\|} \sup_{\mathbf{w} \in \mathbf{W}_0, \mathbf{w} \neq 0} \frac{\|\mathbf{w}\|}{\|\nabla \mathbf{w}\|} \\ &\leq \mathbf{C}_F \frac{\|\nabla \mathbf{q}\|}{\|\mathbf{q}\|}. \end{aligned}$$

Let \mathbf{C}_P be the smallest constant in the inequality

$$\|\mathbf{q}\| \leq \mathbf{C}_P \|\nabla \mathbf{q}\|, \quad \mathbf{q} \in \tilde{\mathbf{W}}(\Omega),$$

i.e.,

$$\inf_{\mathbf{q} \in \tilde{\mathbf{W}}(\Omega), \mathbf{q} \neq 0} \frac{\|\nabla \mathbf{q}\|}{\|\mathbf{q}\|} = \frac{1}{\mathbf{C}_P}.$$

Then

$$\mathbf{C} = \inf_{\mathbf{q} \in \mathring{\mathbf{L}}_2(\Omega), \mathbf{q} \neq 0} \mathcal{E}(\mathbf{q}) \leq \inf_{\mathbf{q} \in \tilde{\mathbf{W}}(\Omega), \mathbf{q} \neq 0} \mathcal{E}(\mathbf{q}) \leq \frac{\mathbf{C}_F}{\mathbf{C}_P}.$$

LBB-condition can be written in the form

$$\|\mathbf{p}\| \leq \mathbf{C}^{-1} \mathbf{I} \nabla \mathbf{p} \mathbf{I} \quad \forall \mathbf{p} \in \mathring{\mathbf{L}}_2(\Omega),$$

what amounts

$$\mathbf{C} \leq \frac{\mathbf{I} \nabla \mathbf{p} \mathbf{I}}{\|\mathbf{p}\|}$$

we see the meaning of this constant: **C is the infimum of \mathbf{H}^{-1} norms of functions such that $\|\mathbf{p}\| = 1$ and $\int_{\Omega} \mathbf{p} \, d\mathbf{x} = 0$.**

Proposition 1

If $\Omega \in \mathbb{R}^n$ then

$$\frac{\|\nabla \mathbf{p}\|_{(-1)}}{\|\mathbf{p}\|} \leq n \quad \forall \mathbf{p} \in \mathbf{L}_2(\Omega).$$

Proof.

$$\sup_{\substack{\mathbf{w} \in \mathbf{W}_0 \\ \mathbf{w} \neq \mathbf{0}}} \frac{\int_{\Omega} \mathbf{p} \operatorname{div} \mathbf{w} \, dx}{\|\nabla \mathbf{w}\|} =$$

$$\sup_{\substack{\mathbf{w} \in \mathbf{W}_0 \\ \mathbf{w} \neq \mathbf{0}}} \frac{\sum_{t=1}^n \int_{\Omega} \mathbf{p} \mathbf{w}_{t,t} \, dx}{\|\nabla \mathbf{w}\|} \leq \sum_{t=1}^n \sup_{\substack{\mathbf{w}_t \in \mathring{\mathbf{H}}^1 \\ \mathbf{w}_t \neq \mathbf{0}}} \frac{\int_{\Omega} \mathbf{p} \mathbf{w}_{t,t} \, dx}{\|\nabla \mathbf{w}_t\|} .$$

Since

$$\|\nabla \mathbf{w}\|^2 = \int_{\Omega} \left(\sum_{t,s=1,n}^n \mathbf{w}_{t,s}^2 \right) dx \geq \int_{\Omega} \mathbf{w}_{t,t}^2 dx \quad \forall t = 1, 2, \dots, n$$

we have

$$\begin{aligned} \sup_{\substack{\mathbf{w} \in \mathbf{W}_0 \\ \mathbf{w} \neq \mathbf{0}}} \frac{\int_{\Omega} \mathbf{p} \operatorname{div} \mathbf{w} dx}{\|\nabla \mathbf{w}\|} &\leq \sum_{t=1}^n \sup_{\substack{\mathbf{w}_t \in \mathring{\mathbf{H}}^1 \\ \mathbf{w}_t \neq \mathbf{0}}} \frac{\int_{\Omega} \mathbf{p} \mathbf{w}_{t,t} dx}{\|\mathbf{w}_{t,t}\|} \leq \\ &\leq \sum_{t=1}^n \sup_{\substack{\eta \in \mathbf{L}_2 \\ \eta \neq \mathbf{0}}} \frac{\int_{\Omega} \mathbf{p} \eta dx}{\|\eta\|} = \sum_{t=1}^n \|\mathbf{p}\| = n \|\mathbf{p}\|. \end{aligned}$$

Proposition 2

If $\mathbf{n} = \mathbf{1}$ then $\mathbf{C} = 1$.

Let $\Omega = (\mathbf{a}, \mathbf{b})$. Due to Proposition 1 we see that $\mathbf{C} \leq 1$. Let \mathbf{p} be an arbitrary function from the set $\mathring{L}_2(\Omega)$. Then, the function

$$\mathbf{w}^{(\mathbf{p})} = \int_{\mathbf{a}}^{\mathbf{x}} \mathbf{p} \, d\mathbf{x} \in \mathbf{W}_0.$$

Really, $\mathbf{w}^{(\mathbf{p})}(\mathbf{a}) = \mathbf{0}$, $\mathbf{w}^{(\mathbf{p})}(\mathbf{b}) = \int_{\mathbf{b}}^{\mathbf{a}} \mathbf{p} \, d\mathbf{x} = \mathbf{0}$ and $\mathbf{w}^{(\mathbf{p})'} = \mathbf{p} \in \mathbf{L}_2(\mathbf{a}, \mathbf{b})$. Thus,

$$\sup_{\substack{\mathbf{w} \in \mathbf{W}_0 \\ \mathbf{w} \neq \mathbf{0}}} \frac{\int_{\Omega} \mathbf{p} \mathbf{w}' \, d\mathbf{x}}{\|\mathbf{w}'\|} \geq \frac{\int_{\Omega} \mathbf{p} \mathbf{w}^{(\mathbf{p})'} \, d\mathbf{x}}{\|\mathbf{w}^{(\mathbf{p})}'\|} = \frac{\int_{\Omega} \mathbf{p}^2 \, d\mathbf{x}}{\|\mathbf{p}\|} = \|\mathbf{p}\|$$

Thus, $\mathbf{C} \geq 1$ and we arrive at the required result.

These estimates give a certain presentation on the value of \mathbf{C} . However, we are mainly interested in the estimate from below, what imposes a task more complicated than the finding the constant in the Friederichs inequality. In principle, one could determine \mathbf{C} by the following arguments. Let $\mathbf{w}_p \in \mathbf{W}_0$ be a function such that

$$\Delta \mathbf{w}_p = \nabla \mathbf{p}, \quad \mathbf{w}_p = \mathbf{0} \text{ on } \partial\Omega.$$

Then,

$$-\int_{\Omega} \nabla \mathbf{w}_p : \nabla \mathbf{v} \, dx = \int_{\Omega} \nabla \mathbf{p} \cdot \mathbf{v} \, dx \quad \forall \mathbf{v} \in \mathbf{W}_0$$

and, thus, we have

$$\int_{\Omega} |\nabla \mathbf{w}_p|^2 \, dx = \int_{\Omega} \mathbf{p} \operatorname{div} \mathbf{w}_p \, dx \quad \forall \mathbf{v} \in \mathbf{W}_0.$$

Therefore,

$$\begin{aligned}
 \mathbf{C} &:= \inf_{\substack{\mathbf{p} \in \mathring{\mathbf{L}}_2(\Omega) \\ \mathbf{p} \neq \mathbf{0}}} \sup_{\substack{\mathbf{w} \in \mathbf{W}_0 \\ \mathbf{w} \neq \mathbf{0}}} \frac{\int_{\Omega} \mathbf{p} \operatorname{div} \mathbf{w} \, dx}{\|\mathbf{p}\| \|\nabla \mathbf{w}\|} \geq \\
 &\geq \inf_{\substack{\mathbf{p} \in \mathring{\mathbf{L}}_2(\Omega) \\ \mathbf{p} \neq \mathbf{0}}} \frac{\int_{\Omega} \mathbf{p} \operatorname{div} \mathbf{w}_p \, dx}{\|\mathbf{p}\| \|\nabla \mathbf{w}_p\|} = \inf_{\substack{\mathbf{p} \in \mathring{\mathbf{L}}_2(\Omega) \\ \mathbf{p} \neq \mathbf{0}}} \frac{\|\nabla \mathbf{w}_p\|}{\|\mathbf{p}\|}.
 \end{aligned}$$

Thus, finding \mathbf{C} requires the minimization of this quotient with respect to all $\mathbf{p} \in \mathring{\mathbf{L}}_2(\Omega)$, where \mathbf{w}_p is taken as the solution of the above defined linear problem. Certainly, such a task (for some Ω) might be solved by only analytical methods.

C for square domains

We will use the above relation to minimize the quotient on a subspace of $\overset{\circ}{L}_2(\Omega)$ what may give a presentation on the value of **C**.

Let

$$\Omega = \mathbb{Q} := \{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{x}_i \in (-\pi, \pi), \mathbf{i} = 1, 2, \dots, n \}.$$

We are interested in the value of the quotient

$$\inf_{\mathbf{p} \in \overset{\circ}{L}_2(\Omega)} \frac{\mathbf{I} \nabla \mathbf{p} \mathbf{I}}{\|\mathbf{p}\|_{\mathbb{Q}}}.$$

Represent **p** as a series with respect to the trial functions

$$\begin{aligned} \mathbf{p}_{ij}^{(1)} &= \sin \mathbf{i} \mathbf{x} \sin \mathbf{j} \mathbf{y}, & \mathbf{p}_{ij}^{(2)} &= \sin \mathbf{i} \mathbf{x} \cos \mathbf{j} \mathbf{y}, \\ \mathbf{p}_{ij}^{(3)} &= \cos \mathbf{i} \mathbf{x} \sin \mathbf{j} \mathbf{y}, & \mathbf{p}_{ij}^{(4)} &= \cos \mathbf{i} \mathbf{x} \cos \mathbf{j} \mathbf{y}, \end{aligned}$$

where $\mathbf{i}, \mathbf{j} = 0, 1, 2, \dots$

Then

$$p(x, y) = \sum_{i,j=0}^{\infty} \sum_{s=1}^4 a_{ij}^{(s)} p_{ij}^{(s)}.$$

Here, the first nonzero coefficients are

$$a_{00}^{(4)} = \frac{1}{4\pi^2} \int_{\Omega} p \, dx dy, \quad a_{i0}^{(2)} = \frac{1}{2\pi^2} \int_{\Omega} p \sin ix \, dx dy,$$

$$a_{0j}^{(3)} = \frac{1}{2\pi^2} \int_{\Omega} p \sin jy \, dx dy,$$

$$a_{i0}^{(4)} = \frac{1}{2\pi^2} \int_{\Omega} p \cos ix \, dx dy,$$

$$a_{0j}^{(2)} = \frac{1}{2\pi^2} \int_{\Omega} p \cos jy \, dx dy,$$

Other coefficients are as follows:

$$\mathbf{a}_{ij}^{(1)} = \frac{1}{\pi^2} \int_{\Omega} \mathbf{p} \sin ix \sin jy \, dx dy ,$$

$$\mathbf{a}_{ij}^{(2)} = \frac{1}{\pi^2} \int_{\Omega} \mathbf{p} \sin ix \cos jy \, dx dy ,$$

$$\mathbf{a}_{ij}^{(3)} = \frac{1}{\pi^2} \int_{\Omega} \mathbf{p} \cos ix \sin jy \, dx dy ,$$

$$\mathbf{a}_{ij}^{(4)} = \frac{1}{\pi^2} \int_{\Omega} \mathbf{p} \cos ix \cos jy \, dx dy .$$

We have

$$\|\mathbf{p}\|_{\mathbb{Q}}^2 = \pi^2 \sum_{i,j=0}^{\infty} \lambda_{ij} \left[\left(\mathbf{a}_{ij}^{(1)}\right)^2 + \left(\mathbf{a}_{ij}^{(2)}\right)^2 + \left(\mathbf{a}_{ij}^{(3)}\right)^2 + \left(\mathbf{a}_{ij}^{(4)}\right)^2 \right] ,$$

where $\lambda_{00} = 0$, $\lambda_{01} = 2$, $\lambda_{10} = 2$ and $\lambda_{ij} = 1$ for all $\mathbf{i}, \mathbf{j} \geq \mathbf{1}$.

Let us take a finite number of elements in the Fourier series for \mathbf{p} :

$$\mathbf{p} = \sum_{i,j=0}^N \sum_{s=1}^4 \mathbf{a}_{ij}^{(s)} \mathbf{p}_{ij}^{(s)},$$

where $\mathbf{a}_{ij}^{(s)}$ are the above defined coefficients. Since

$$\mathbf{I} \nabla \mathbf{p} \mathbf{I} = \sup_{\mathbf{v} \in \mathbb{W}_0} \frac{\int_{\Omega} \mathbf{p} \operatorname{div} \mathbf{v} \, dx}{\|\nabla \mathbf{v}\|_{\mathbb{Q}}}$$

we need to introduce the system of trial functions in $\mathbb{W}_0(\mathbb{Q})$. It is given by the system of eigenfunctions for the problem

$$\Delta \mathbf{w} = \mu \mathbf{w} \quad \mathbf{w}|_{\partial \mathbb{Q}} = \mathbf{0}.$$

This system is

$$\phi_{\alpha\beta} = \sin \frac{\alpha}{2}(x + \pi) \sin \frac{\beta}{2}(y + \pi).$$

In this case,

$$\phi_{\alpha\beta,1} = \frac{\alpha}{2} \cos \frac{\alpha}{2}(x + \pi) \sin \frac{\beta}{2}(y + \pi),$$

$$\phi_{\alpha\beta,2} = \frac{\beta}{2} \sin \frac{\alpha}{2}(x + \pi) \cos \frac{\beta}{2}(y + \pi).$$

Take a finite number \mathbf{M} of basic functions in the representation of \mathbf{v} , namely we set

$$\mathbf{v} = \mathbf{v}^{\mathbf{M}} = (\mathbf{v}_1^{\mathbf{M}}, \mathbf{v}_2^{\mathbf{M}}), \quad \mathbf{v}_1^{\mathbf{M}} = \sum_{\alpha,\beta=1}^{\mathbf{M}} \mathbf{b}_{\alpha\beta} \phi_{\alpha\beta}, \quad \mathbf{v}_2^{\mathbf{M}} = \sum_{\alpha,\beta=1}^{\mathbf{M}} \mathbf{c}_{\alpha\beta} \phi_{\alpha\beta}.$$

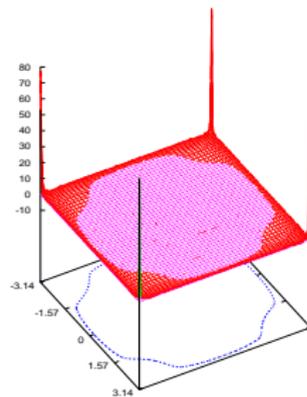
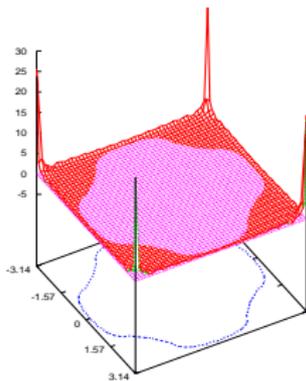
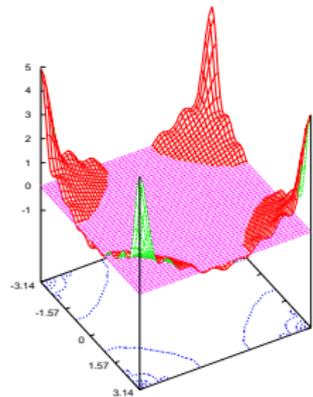
The set of all such functions we denote $\mathbf{W}_0^{\mathbf{M}}$. In this case, we can obtain a lower bound for the required norm. Really, we have

$$\mathbf{I} \nabla \mathbf{p} \mathbf{I}^{(\mathbf{M})} := \sup_{\mathbf{v}^{\mathbf{M}} \in \mathbf{W}_0^{\mathbf{M}}} \frac{\int_{\Omega} \mathbf{p} \operatorname{div} \mathbf{v}^{\mathbf{M}} \, dx}{\|\nabla \mathbf{v}^{\mathbf{M}}\|_{\mathbb{Q}}} \leq \mathbf{I} \nabla \mathbf{p} \mathbf{I} = \sup_{\mathbf{v} \in \mathbf{W}_0} \frac{\int_{\Omega} \mathbf{p} \operatorname{div} \mathbf{v} \, dx}{\|\nabla \mathbf{v}\|_{\mathbb{Q}}}.$$

Thus, we may hope to estimate the value of the quotient

$$\inf_{\mathbf{p} \in \overset{\circ}{\mathbf{L}}_2(\Omega)} \frac{\mathbf{I} \nabla \mathbf{p} \mathbf{I}}{\|\mathbf{p}\|_{\mathbb{Q}}}.$$

by taking $\mathbf{N}, \mathbf{M} \rightarrow +\infty$, $\mathbf{M} = \kappa \mathbf{N}$ κ is essentially larger than 1 (typically 8-20). Numerical results for different \mathbf{N} are exposed below.



Minimizer p_n for $n = 8, 36$ and 120 .

Deviation estimates for the Stokes problem

In order to clarify the main ideas of our approach we rewrite the classical Stokes system in a somewhat different form:

$$\mathbf{div} \boldsymbol{\sigma} = \nabla \mathbf{p} - \mathbf{f} \quad \text{in } \Omega, \quad (253)$$

$$\mathbf{div} \mathbf{u} = \mathbf{0} \quad \text{in } \Omega, \quad (254)$$

$$\boldsymbol{\sigma} = \nu \nabla \mathbf{u} \quad \text{in } \Omega, \quad (255)$$

$$\mathbf{u} = \mathbf{0} \quad \text{on } \partial\Omega. \quad (256)$$

This system involves one additional variable $\boldsymbol{\sigma}$ that corresponds to the field of stresses. Now we may regard the Stokes problem as the problem of finding a triplet of functions $(\mathbf{u}, \boldsymbol{\sigma}, \mathbf{p})$.

Primal and Dual Problems

Functional formulations of the above problem are given in natural "energy" set for this velocity–stress–pressure setting, which is

$$\mathcal{E} := \mathring{\mathbf{J}}_2^1(\Omega) \times \Sigma \times \mathring{\mathbf{L}}_2(\Omega).$$

Problem \mathcal{P} . Find $\mathbf{u} \in \mathring{\mathbf{J}}_2^1(\Omega)$ such that

$$\mathbf{J}(\mathbf{u}) \leq \mathbf{J}(\mathbf{v}) \quad \text{for all } \mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega),$$

where

$$\mathbf{J}(\mathbf{v}) = \int_{\Omega} \left(\frac{\nu}{2} |\nabla \mathbf{v}|^2 - \mathbf{f} \cdot \mathbf{v} \right) \mathbf{d}\mathbf{x}.$$

We denote the exact lower bound of this problem by $\inf \mathcal{P}$.

Introduce the Lagrangian $\mathbf{L} : \mathring{\mathbf{J}}_2^1(\Omega) \times \Sigma(\Omega) \rightarrow \mathbb{R}$:

$$\mathbf{L}(\mathbf{v}, \boldsymbol{\tau}) = \int_{\Omega} \left(\boldsymbol{\tau} : \nabla \mathbf{v} - \frac{1}{2\nu} |\boldsymbol{\tau}|^2 \right) \mathbf{d}\mathbf{x} - \int_{\Omega} \mathbf{f}\mathbf{v} \mathbf{d}\mathbf{x}$$

that together with Problem \mathcal{P} generates the dual problem

$$\sup_{\boldsymbol{\tau} \in \Sigma} \inf_{\mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega)} \mathbf{L}(\mathbf{v}, \boldsymbol{\tau})$$

which is **Problem \mathcal{P}^*** : find $\boldsymbol{\sigma} \in \Sigma_{\mathbf{f}}$ such that

$$\mathbf{I}^*(\boldsymbol{\sigma}) = \sup_{\boldsymbol{\tau} \in \Sigma_{\mathbf{f}}} \mathbf{I}^*(\boldsymbol{\tau}), \quad \mathbf{I}^*(\boldsymbol{\tau}) = -\frac{1}{2\nu} \int_{\Omega} |\boldsymbol{\tau}|^2 \mathbf{d}\mathbf{x}$$

where

$$\Sigma_{\mathbf{f}} := \left\{ \boldsymbol{\tau} \in \Sigma(\Omega) \mid \int_{\Omega} \boldsymbol{\tau} : \nabla \mathbf{w} \mathbf{d}\mathbf{x} = \int_{\Omega} \mathbf{f}\mathbf{w} \mathbf{d}\mathbf{x} \text{ for all } \mathbf{w} \in \mathring{\mathbf{J}}_2^1(\Omega) \right\}.$$

From the general theorems of convex analysis it follows

Theorem (1)

There exists a unique minimizer \mathbf{u} of problem \mathcal{P} and unique maximizer $\boldsymbol{\sigma}$ of problem \mathcal{P}^ . These two functions meet the equalities*

$$\mathbf{I}^*(\boldsymbol{\sigma}) = \sup \mathcal{P}^* = \inf \mathcal{P} = \mathbf{I}(\mathbf{u}), \quad (257)$$

$$\boldsymbol{\sigma} = \nu \nabla \mathbf{u}. \quad (258)$$

Basic error estimate

The basic error relation for the Stokes problem is given by the following theorem (S. Repin, 2002).

Theorem (2)

For any $\mathbf{v} \in \mathbf{J}_2^1(\Omega)$ and any $\tau_f \in \Sigma_f$, we have

$$\nu \|\nabla(\mathbf{v} - \mathbf{u})\|^2 + \frac{1}{\nu} \|\tau_f - \sigma\|^2 = 2(\mathbf{J}(\mathbf{v}) - \mathbf{I}^*(\tau_f)). \quad (259)$$

Proof of Theorem 2

Since \mathbf{u} is the solution of problem \mathcal{P} , we obtain

$$\begin{aligned} \mathbf{J}(\mathbf{v}) - \mathbf{J}(\mathbf{u}) &= \int_{\Omega} \left(\frac{\nu}{2} |\nabla \mathbf{v}|^2 - \frac{\nu}{2} |\nabla \mathbf{u}|^2 - \mathbf{f} \cdot (\mathbf{v} - \mathbf{u}) \right) \mathbf{d}\mathbf{x} = \\ &= \int_{\Omega} \left(\frac{\nu}{2} |\nabla(\mathbf{v} - \mathbf{u})|^2 + \nu \nabla \mathbf{u} : \nabla(\mathbf{v} - \mathbf{u}) - \mathbf{f} \cdot (\mathbf{v} - \mathbf{u}) \right) \mathbf{d}\mathbf{x} = \\ &= \frac{\nu}{2} \int_{\Omega} |\nabla(\mathbf{v} - \mathbf{u})|^2 \mathbf{d}\mathbf{x} \quad \text{for all } \mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega). \end{aligned}$$

Since $\mathbf{J}(\mathbf{u}) = \inf \mathcal{P}$, we conclude that

$$\frac{\nu}{2} \|\nabla(\mathbf{v} - \mathbf{u})\|^2 = \mathbf{J}(\mathbf{v}) - \inf \mathcal{P} \quad \text{for all } \mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega).$$

The next step is to derive a similar relation for the dual problem. For this purpose, we note that the maximizer σ of problem \mathcal{P}^* satisfies the relation

$$\int_{\Omega} \sigma : (\tau_f - \sigma) \, \mathbf{d}\mathbf{x} = \mathbf{0} \quad \text{for all } \tau_f \in \Sigma_f.$$

By virtue of this relation, we find that

$$\sup \mathcal{P}^* - \mathbf{I}^*(\tau_f) = \mathbf{I}^*(\sigma) - \mathbf{I}^*(\tau_f) = \frac{1}{2\nu} \|\tau_f - \sigma\|^2 \quad \tau_f \in \Sigma_f.$$

Since $\inf \mathcal{P} = \sup \mathcal{P}^*$ we sum the two equalities and obtain

$$\nu \|\nabla(\mathbf{v} - \mathbf{u})\|^2 + \frac{1}{\nu} \|\tau_f - \sigma\|^2 = 2(\mathbf{J}(\mathbf{v}) - \mathbf{I}^*(\tau_f)).$$

Stokes problem is a particular case of the abstract problem we investigated in Lecture 5:

Find $\mathbf{u} \in \mathbf{V}_0 + \mathbf{u}_0$ such that

$$(\mathcal{A}\mathbf{\Lambda}\mathbf{u}, \mathbf{\Lambda}\mathbf{w}) + \langle \ell, \mathbf{w} \rangle = 0 \quad \forall \mathbf{w} \in \mathbf{V}_0.$$

In this case $\mathbf{V}_0 = \mathbf{J}_2^1(\Omega)$, \mathbf{V} is a subspace of \mathbf{H}^1 containing solenoidal fields, $\mathbf{\Lambda} = \nabla$ (tensor-gradient), $\mathbf{U} = \mathbf{\Sigma}$, $\mathcal{A}\mathbf{y} = \nu\mathbf{y}$, and

$$\langle \ell, \mathbf{w} \rangle = - \int_{\Omega} \mathbf{f}\mathbf{w} \, dx$$

Thus, we can apply the estimate

$$\frac{1}{2} \|\Lambda(\mathbf{v} - \mathbf{u})\|^2 \leq (1 + \beta) \mathbf{D}(\Lambda \mathbf{v}, \mathbf{y}) + \frac{1 + \beta}{2\beta} \mathbf{I} \ell + \Lambda^* \mathbf{y} \mathbf{I}^2, \quad (260)$$

where $\|\mathbf{y}\|^2 = \int_{\Omega} \nu |\mathbf{y}|^2 \mathbf{d}\mathbf{x}$ and

$$\begin{aligned} \mathbf{I} \ell + \Lambda^* \mathbf{y} \mathbf{I} &= \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\langle \ell + \Lambda^* \mathbf{y}, \mathbf{w} \rangle}{\|\Lambda \mathbf{w}\|} = \sup_{\mathbf{w} \in \mathring{\mathbf{H}}_2(\Omega)} \frac{\int_{\Omega} (\nabla \mathbf{w} : \mathbf{y} - \mathbf{f}\mathbf{w}) \mathbf{d}\mathbf{x}}{\|\nabla \mathbf{w}\|} = \\ &\sup_{\mathbf{w} \in \mathring{\mathbf{H}}_2(\Omega)} \frac{\int_{\Omega} (\nabla \mathbf{w} : \mathbf{y} - \mathbf{f}\mathbf{w} - \mathbf{q} \operatorname{div} \mathbf{w}) \mathbf{d}\mathbf{x}}{\|\nabla \mathbf{w}\|} \leq \\ &\leq \sup_{\mathbf{w} \in \mathring{\mathbf{H}}^1(\Omega)} \frac{\int_{\Omega} (\nabla \mathbf{w} : \mathbf{y} - \mathbf{f}\mathbf{w} - \mathbf{q} \operatorname{div} \mathbf{w}) \mathbf{d}\mathbf{x}}{\|\nabla \mathbf{w}\|} \quad \forall \mathbf{q} \in \mathbf{L}^2(\Omega). \end{aligned}$$

If

$$\mathbf{y} \in \Sigma_{\text{div}}(\Omega) := \{\mathbf{y} \in \Sigma \mid \text{div} \mathbf{y} \in \mathbf{L}^2(\Omega, \mathbb{R}^n)\}$$

and $\mathbf{q} \in \mathbf{H}^1$, we have

$$\sup_{\mathbf{w} \in \mathring{\mathbf{H}}^1(\Omega)} \frac{\int_{\Omega} (\nabla \mathbf{w} : \mathbf{y} - \mathbf{f} \mathbf{w} - \mathbf{q} \text{div} \mathbf{w}) \, dx}{\|\nabla \mathbf{w}\|} = \sup_{\mathbf{w} \in \mathring{\mathbf{H}}^1(\Omega)} \frac{\int_{\Omega} (\mathbf{f} - \nabla \mathbf{q} + \text{div} \mathbf{y}) \cdot \mathbf{w} \, dx}{\|\nabla \mathbf{w}\|}$$

Since

$$\|\mathbf{w}\| \leq \mathbf{C}_{\Omega} \|\nabla \mathbf{w}\| = \mathbf{C}_{\Omega} \nu^{-1/2} \|\nabla \mathbf{w}\|,$$

we obtain

$$|\ell + \Lambda^* \mathbf{y}| \leq \mathbf{C}_{\Omega} \nu^{-1/2} \|\mathbf{f} - \nabla \mathbf{q} + \text{div} \mathbf{y}\|$$

Further,

$$\begin{aligned} \mathbf{D}(\nabla \mathbf{v}, \mathbf{y}) &= \int_{\Omega} \left(\frac{1}{2} \nu \nabla \mathbf{v} : \nabla \mathbf{v} + \frac{1}{2} \nu^{-1} \mathbf{y} : \mathbf{y} - \nabla \mathbf{v} : \mathbf{y} \right) \mathbf{d}\mathbf{x} = \\ &= \frac{1}{2\nu} \|\mathbf{y} - \nu \nabla \mathbf{v}\|^2. \end{aligned}$$

Now, from (260) we obtain

$$\frac{\nu}{2} \|\nabla(\mathbf{u} - \mathbf{v})\|^2 \leq (1 + \beta) \frac{1}{2\nu} \|\mathbf{y} - \nu \nabla \mathbf{v}\|^2 + \frac{1 + \beta}{2\beta\nu} \mathbf{C}_{\Omega}^2 \|\mathbf{f} - \nabla \mathbf{q} + \mathbf{div} \mathbf{y}\|^2,$$

or

$$\nu^2 \|\nabla(\mathbf{u} - \mathbf{v})\|^2 \leq (1 + \beta) \|\mathbf{y} - \nu \nabla \mathbf{v}\|^2 + \frac{1 + \beta}{\beta} \mathbf{C}_{\Omega}^2 \|\mathbf{f} - \nabla \mathbf{q} + \mathbf{div} \mathbf{y}\|^2.$$

Deviation estimate for solenoidal approximations

By the minimization with respect to β we derive the first basic estimate for the Stokes problem:

$$\nu \|\nabla(u - \mathbf{v})\| \leq \|\mathbf{y} - \nu \nabla \mathbf{v}\| + \mathbf{C}_\Omega \|\mathbf{f} - \nabla \mathbf{q} + \operatorname{div} \mathbf{y}\|. \quad (261)$$

Here \mathbf{v} is **any** conforming approximation of \mathbf{u} and \mathbf{y} is **any** tensor-function in $\Sigma_{\operatorname{div}}(\Omega)$ and $\mathbf{q} \in \mathbf{H}^1$ is an "image" of the pressure function.

This and the next estimate for non-solenoidal approximations has been derived in '99, English translation is presented in S. Repin. A posteriori estimates for the Stokes problem, *J. Math. Sci. (New York)*, **109** (2002).

Non-solenoidal approximations

If the function $\widehat{\mathbf{v}} \in \mathbf{V}_0 + \mathbf{u}_0$ does not satisfy the incompressibility condition, then the estimate of its deviation from \mathbf{u} can be obtained as follows.

By Lemma 2 for the function $\widehat{\mathbf{v}}_0 := \widehat{\mathbf{v}} - \mathbf{u}_0$ one can find a function $\mathbf{w}_0 \in \mathring{\mathbf{J}}_2^1(\Omega)$ such that

$$\|\nabla(\widehat{\mathbf{v}}_0 - \widehat{\mathbf{w}}_0)\| \leq \kappa_2(\Omega) \|\operatorname{div} \widehat{\mathbf{v}}_0\|.$$

Then,

$$\begin{aligned} \nu \|\nabla(\mathbf{u} - \widehat{\mathbf{v}})\| &= \nu \|\nabla(\mathbf{u} - \widehat{\mathbf{v}}_0 - \mathbf{u}_0)\| \leq \\ &\leq \nu \|\nabla(\mathbf{u} - (\widehat{\mathbf{w}}_0 + \mathbf{u}_0))\| + \nu \|\nabla(\widehat{\mathbf{v}}_0 - \widehat{\mathbf{w}}_0)\|. \end{aligned}$$

Use (261) to estimate the first norm in the right-hand side of this inequality.

We obtain

$$\begin{aligned} \nu \|\nabla(\mathbf{u} - \hat{\mathbf{v}})\| &\leq \|\nu \nabla(\hat{\mathbf{w}}_0 + \mathbf{u}_0) - \mathbf{y}\| + \mathbf{C}_\Omega \|\mathbf{div} \mathbf{y} + \mathbf{f} - \nabla \mathbf{q}\| + \\ &+ \nu \|\nabla(\hat{\mathbf{v}}_0 - \hat{\mathbf{w}}_0)\| \leq \|\nu \nabla \hat{\mathbf{v}} - \mathbf{y}\| + \\ &+ \mathbf{C}_\Omega \|\mathbf{div} \mathbf{y} + \mathbf{f} - \nabla \mathbf{q}\| + 2\nu \|\nabla(\hat{\mathbf{v}}_0 - \hat{\mathbf{w}}_0)\|. \end{aligned}$$

Hence, we arrive at the estimate

$$\boxed{\nu \|\nabla(\mathbf{u} - \hat{\mathbf{v}})\| \leq \|\nu \nabla(\hat{\mathbf{v}}) - \mathbf{y}\| + \mathbf{C}_\Omega \|\mathbf{div} \mathbf{y} + \mathbf{f} - \nabla \mathbf{q}\| + \frac{2\nu}{\mathbf{C}} \|\mathbf{div} \hat{\mathbf{v}}\|.} \quad (262)$$

Three terms in the right-hand side of the estimate present three natural parts of the error, namely **errors in the constitutive law, differential equation and incompressibility condition**.

Another form of the Majorant

Set $\mathbf{y} = \boldsymbol{\eta} + \mathbf{q}\mathbb{I}$, where \mathbb{I} is the unit tensor and $\boldsymbol{\eta} \in \boldsymbol{\Sigma}_{\text{div}}(\boldsymbol{\Omega})(\boldsymbol{\Omega})$. Then the Majorant comes in the form

$$\nu \|\nabla(\mathbf{u} - \widehat{\mathbf{v}})\| \leq \|\nu \nabla(\widehat{\mathbf{v}}) - \boldsymbol{\eta} - \mathbf{q}\mathbb{I}\| + \mathbf{C}_{\Omega} \|\text{div} \boldsymbol{\eta} + \mathbf{f}\| + \frac{2\nu}{\mathbf{C}} \|\text{div} \widehat{\mathbf{v}}\|. \quad (263)$$

Thus, if the constants \mathbf{c}_{Ω} and \mathbf{C} are known (or we know suitable upper bounds for them), then (262) and (263) provides a way of practical estimation the deviation of $\widehat{\mathbf{v}}$ from \mathbf{u} .

Practical implementation

To use the above estimates in practice we should select certain finite dimensional subspaces

$$\Sigma_k \quad \text{and} \quad Q_k$$

for the functions \mathbf{y} (or $\boldsymbol{\eta}$) and \mathbf{q} , respectively.

Minimization of the right-hand side of the estimates with respect to \mathbf{y} and \mathbf{q} gives an estimate of the deviation, which will be the sharper the greater is the dimensionality of the subspaces used.

Estimates for the pressure field

Let $\mathbf{q} \in \mathring{\mathbf{L}}_2(\Omega)$ be an approximation of the pressure field \mathbf{p} . Then $(\mathbf{p} - \mathbf{q}) \in \mathring{\mathbf{L}}_2(\Omega)$ and the Inf-Sup condition implies the relation

$$\sup_{\mathbf{w} \in \mathbf{V}_0, \mathbf{w} \neq \mathbf{0}} \frac{\int_{\Omega} (\mathbf{p} - \mathbf{q}) \operatorname{div} \mathbf{w} \, d\mathbf{x}}{\|\mathbf{p} - \mathbf{q}\| \|\nabla \mathbf{w}\|} \geq \mathbf{C}.$$

Thus, for any small positive ϵ there exists a nonzero function $\mathbf{w}_{\mathbf{p}\mathbf{q}}^{\epsilon} \in \mathbf{V}_0$ such that

$$\int_{\Omega} (\mathbf{p} - \mathbf{q}) \operatorname{div} \mathbf{w}_{\mathbf{p}\mathbf{q}}^{\epsilon} \, d\mathbf{x} \geq (\mathbf{C} - \epsilon) \|\mathbf{p} - \mathbf{q}\| \|\nabla \mathbf{w}_{\mathbf{p}\mathbf{q}}^{\epsilon}\|.$$

Since

$$\int_{\Omega} \nu \varepsilon(\mathbf{u}) : \varepsilon(\mathbf{w}_{pq}^{\varepsilon}) \, d\mathbf{x} = \int_{\Omega} (\mathbf{f} \cdot \mathbf{w}_{pq}^{\varepsilon} + \mathbf{p} \operatorname{div} \mathbf{w}_{pq}^{\varepsilon}) \, d\mathbf{x},$$

we have

$$\begin{aligned} & \int_{\Omega} (\mathbf{p} - \mathbf{q}) \operatorname{div} \mathbf{w}_{pq}^{\varepsilon} \, d\mathbf{x} = \\ & = \int_{\Omega} \left\{ \nu \varepsilon((\mathbf{u} - \widehat{\mathbf{v}})) : \varepsilon(\mathbf{w}_{pq}^{\varepsilon}) + (\nu \varepsilon(\widehat{\mathbf{v}})) : \varepsilon(\mathbf{w}_{pq}^{\varepsilon}) + \nabla \mathbf{q} \cdot \mathbf{w}_{pq}^{\varepsilon} - \mathbf{f} \cdot \mathbf{w}_{pq}^{\varepsilon} \right\} \, d\mathbf{x} \\ & = \int_{\Omega} \nu \varepsilon((\mathbf{u} - \widehat{\mathbf{v}})) : \varepsilon(\mathbf{w}_{pq}^{\varepsilon}) \, d\mathbf{x} + \int_{\Omega} (\nu \varepsilon(\widehat{\mathbf{v}} - \mathbf{y})) : \varepsilon(\mathbf{w}_{pq}^{\varepsilon}) \, d\mathbf{x} \\ & \quad + \int_{\Omega} (\mathbf{y} : \varepsilon(\mathbf{w}_{pq}^{\varepsilon}) + \nabla \mathbf{q} \cdot \mathbf{w}_{pq}^{\varepsilon} - \mathbf{f} \cdot \mathbf{w}_{pq}^{\varepsilon}) \, d\mathbf{x}, \end{aligned}$$

where $\widehat{\mathbf{v}}$ is an arbitrary function in $\mathbf{W}_0 + \mathbf{u}_0$ and \mathbf{y} as an arbitrary tensor-valued function in Σ .

Above relations lead to the estimates

$$\begin{aligned} \|\mathbf{p} - \mathbf{q}\| &\leq \frac{1}{(\mathbf{C} - \epsilon)\|\nabla \mathbf{w}_{\mathbf{p}\mathbf{q}}^\epsilon\|} \\ &\times \left[\int_{\Omega} (\nu \varepsilon(\mathbf{u} - \widehat{\mathbf{v}}) : \varepsilon(\mathbf{w}_{\mathbf{p}\mathbf{q}}^\epsilon) + (\nu \varepsilon(\widehat{\mathbf{v}}) - \mathbf{y}) : \varepsilon(\mathbf{w}_{\mathbf{p}\mathbf{q}}^\epsilon)) \, \mathbf{d}\mathbf{x} \right. \\ &\quad \left. + \int_{\Omega} (-\mathbf{w}_{\mathbf{p}\mathbf{q}}^\epsilon \cdot \mathbf{div} \mathbf{y} + \nabla \mathbf{q} \cdot \mathbf{w}_{\mathbf{p}\mathbf{q}}^\epsilon - \mathbf{f} \cdot \mathbf{w}_{\mathbf{p}\mathbf{q}}^\epsilon) \, \mathbf{d}\mathbf{x} \right] \\ &\leq \frac{1}{(\mathbf{C} - \epsilon)} [\nu \|\varepsilon(\mathbf{u} - \widehat{\mathbf{v}})\| + \|\nu \varepsilon(\widehat{\mathbf{v}}) - \mathbf{y}\| + \mathbf{C}_\Omega \|\mathbf{div} \mathbf{y} + \mathbf{f} - \nabla \mathbf{q}\|]. \end{aligned}$$

The first term in the right-hand side of this inequality is estimated by (262).

Deviation estimate for the pressure function

Since ϵ may be taken arbitrarily small, we obtain the following estimate for the deviation from the exact pressure field:

$$\begin{aligned} \frac{1}{2} \|\mathbf{p} - \mathbf{q}\| &\leq \frac{\nu}{\mathbf{C}^2} \|\mathbf{div} \hat{\mathbf{v}}\| + \\ &+ \frac{1}{\mathbf{C}} \|\nu \epsilon(\hat{\mathbf{v}}) - \mathbf{y}\| + \frac{\mathbf{C}_\Omega}{\mathbf{C}} \|\mathbf{div} \mathbf{y} + \mathbf{f} - \nabla \mathbf{q}\|. \end{aligned} \quad (264)$$

It is easy to see that the right-hand side of (264) consists of the same terms as the right-hand side of (262) and vanishes if and only if, $\hat{\mathbf{v}} = \mathbf{u}$, $\mathbf{y} = \boldsymbol{\sigma}$ and $\mathbf{p} = \mathbf{q}$. However, in this case, the dependence of the penalty multipliers from the constant \mathbf{C} is stronger.

Problems with condition $\operatorname{div} \mathbf{u} = \phi$.

In many cases, divergence-free condition is replaced by

$$\operatorname{div} \mathbf{u} = \phi \quad \text{in } \Omega,$$

where ϕ is a given function in $\mathring{L}_2(\Omega)$. For such functions, we have the problem: find \mathbf{u} that is equal to \mathbf{u}_0 on $\partial\Omega$ and

$$\begin{aligned} -\operatorname{div} \boldsymbol{\sigma} + \nabla \mathbf{p} &= \mathbf{f} \quad \text{in } \Omega, \\ \boldsymbol{\sigma} &= \nu \boldsymbol{\varepsilon}(\mathbf{u}) \quad \text{in } \Omega, \end{aligned}$$

Let $\mathbf{u}_\phi \in \mathbf{W}_0$, $\operatorname{div} \mathbf{u}_\phi = \phi$. By setting $\mathbf{u} = \bar{\mathbf{u}} + \mathbf{u}_\phi$ and $\bar{\mathbf{u}}_0 = \mathbf{u}_0 - \mathbf{u}_\phi$, we present the boundary-value problem as follows: find $\bar{\mathbf{u}} \in \mathring{\mathbf{J}}_2^1(\Omega) + \bar{\mathbf{u}}_0$ such that

$$\begin{aligned} -\operatorname{div} \bar{\boldsymbol{\sigma}} + \nabla \mathbf{p} &= \bar{\mathbf{f}} \quad \text{in } \Omega, & \bar{\mathbf{f}} &= \mathbf{f} + \nu \operatorname{div} \boldsymbol{\varepsilon}(\mathbf{u}_\phi) \in \mathbf{H}^{-1}, \\ \bar{\boldsymbol{\sigma}} &= \nu \boldsymbol{\varepsilon}(\bar{\mathbf{u}}) \quad \text{in } \Omega. \end{aligned}$$

Assume that \mathbf{u} is approximated by a certain $\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0$. Let \mathbf{v} be presented in the form $\mathbf{v} = \bar{\mathbf{v}} + \mathbf{u}_\phi$. Now, we apply (262) to a "shifted" system and obtain

$$\begin{aligned} \|\varepsilon(\mathbf{u} - \mathbf{v})\| &= \|\varepsilon(\bar{\mathbf{u}} - \bar{\mathbf{v}})\| \leq \\ &\leq \|\nu\varepsilon(\bar{\mathbf{v}}) - \mathbf{y}\| + \mathbf{I} \mathbf{divy} + \bar{\mathbf{f}} - \nabla\mathbf{q} \mathbf{I} + \frac{2\nu}{\mathbf{C}_{\text{LBB}}} \|\mathbf{div}\bar{\mathbf{v}}\|. \end{aligned}$$

Set here $\mathbf{y} = -\nu\nabla(\mathbf{u}_\phi) + \boldsymbol{\eta}$, where $\boldsymbol{\eta}$ is a function in Σ_{div} . Then

$$\mathbf{divy} + \bar{\mathbf{f}} = -\nu\mathbf{div}\varepsilon(\mathbf{u}_\phi) + \mathbf{div}\boldsymbol{\eta} + \bar{\mathbf{f}} = \mathbf{div}\boldsymbol{\eta} + \mathbf{f}$$

and $\nu\varepsilon(\bar{\mathbf{v}}) - \mathbf{y} = \nu\varepsilon(\mathbf{v} - \mathbf{u}_\phi) - \mathbf{y} = \nu\varepsilon(\mathbf{v}) - \boldsymbol{\eta}$. Therefore,

$$\begin{aligned} \|\varepsilon(\mathbf{u} - \mathbf{v})\| &\leq \\ &\leq \|\nu\varepsilon(\mathbf{v}) - \boldsymbol{\eta}\| + \mathbf{I} \mathbf{div}\boldsymbol{\eta} + \mathbf{f} - \nabla\mathbf{q} \mathbf{I} + \frac{2\nu}{\mathbf{C}_{\text{LBB}}} \|\mathbf{div}\mathbf{v} - \phi\|. \end{aligned}$$

Problems for almost incompressible fluids

Models of almost incompressible fluids are often used for constructing sequences of functions converging to a solution of the Stokes problem. In this case, the incompressibility condition is replaced by a penalty term: find $\mathbf{u}_\delta \in \mathbf{V}$ satisfying the integral identity

$$\int_{\Omega} (\nu \nabla \mathbf{u}_\delta : \nabla \mathbf{w} + \frac{1}{\delta} \operatorname{div} \mathbf{u}_\delta \operatorname{div} \mathbf{w}) \, dx = \int_{\Omega} \mathbf{f} \cdot \mathbf{w} \, dx, \quad \mathbf{w} \in \mathbf{W}_0,$$

and the boundary condition $\mathbf{u}_\delta = \mathbf{u}_0$ on $\partial\Omega$.

It is not difficult to show (see, e.g., [R. Temam. NS Equations](#)), that \mathbf{u}_δ tends to \mathbf{u} (solution of the Stokes problem) in \mathbf{H}^1 norm and

$\mathbf{p}_\delta = -\frac{1}{\delta} \operatorname{div} \mathbf{u}_\delta \in \mathring{\mathbf{L}}_2(\Omega)$ converges to the respective pressure function \mathbf{p} in \mathbf{L}_2 as $\delta \rightarrow 0$.

By (262) we can easily obtain an estimate of the difference between \mathbf{u} and \mathbf{u}_δ . Let us set in (262) $\mathbf{y} = \boldsymbol{\tau}_\delta := \nu \nabla \mathbf{u}_\delta$ and $\mathbf{q} = \mathbf{p}_\delta = -\frac{1}{\delta} \operatorname{div} \mathbf{u}_\delta$. In this case, $\|\nu \nabla \mathbf{u}_\delta - \boldsymbol{\tau}_\delta\| = \mathbf{0}$ and

$$\begin{aligned} \mathbf{[div} \boldsymbol{\tau}_\delta + \mathbf{f} - \nabla \mathbf{p}_\delta \mathbf{]} &= \\ &= \sup_{\mathbf{w} \in \mathbf{V}_0} \frac{\int_{\Omega} (-\nu \nabla \mathbf{u}_\delta : \nabla \mathbf{w} + \mathbf{f} \cdot \mathbf{w} + \mathbf{p}_\delta \operatorname{div} \mathbf{w}) \, \mathbf{d}\mathbf{x}}{\|\nabla \mathbf{w}\|} = \mathbf{0}. \end{aligned}$$

Thus, we conclude that

$$\frac{1}{2} \|\nabla(\mathbf{u} - \mathbf{u}_\delta)\| \leq \frac{1}{\mathbf{C}_{\text{LBB}}} \|\operatorname{div} \mathbf{u}_\delta\|,$$

We observe that the deviation from the exact solution of the Stokes problem is controlled by the norm of the divergence of the regularized problem. Similar estimate can be obtained for the approximations constructed by means of the Uzawa algorithm.

Functional a posteriori estimates for the Stokes and some other problems were also derived by **nonvariational** techniques (see [S. Repin. *St.-Petersburg Math. J.*, \(2004\)](#)).

In particular, such estimates were derived for the **Oseen problem**

$$\begin{aligned} -\nu \Delta \mathbf{u} + \operatorname{div}(\mathbf{a} \otimes \mathbf{u}) &= \mathbf{f} - \nabla p && \text{in } \Omega, \\ \operatorname{div} \mathbf{u} &= 0 && \text{in } \Omega, \\ \mathbf{u} &= \mathbf{0} && \text{on } \partial\Omega. \end{aligned}$$

Generalizations

A posteriori estimates of the above discussed type can be derived in the abstract form for the whole class of problems where a solution is seeking in a subspace.

Typically, we have the following diagram:

$$\begin{array}{ccccc}
 \mathbf{H} & \xleftarrow{\mathbf{B}} & \mathbf{W}_0 & \xrightarrow{\Lambda} & \mathbf{U} \quad (\mathbf{Y}, \mathbf{Y}^*) \\
 & & \updownarrow & & \\
 \mathbf{H} & \xrightarrow{\mathbf{B}^*} & \mathbf{W}_0^* & \xleftarrow{\Lambda^*} & \mathbf{U}
 \end{array}$$

Basic problem. Find $\mathbf{p} \in \mathbf{H}$ and $\mathbf{u} \in \mathcal{V}_0$ that satisfy the relation

$$(\mathcal{A}\Lambda\mathbf{u}, \Lambda\mathbf{w}) + \langle \mathbf{f} - \mathbf{B}^*\mathbf{p}, \mathbf{w} \rangle = 0 \quad \forall \mathbf{w} \in \mathbf{W}_0,$$

where

$$\mathcal{V}_0 = \text{Ker}\mathbf{B} := \{\mathbf{v} \in \mathbf{W}_0 \mid \mathbf{B}\mathbf{v} = \mathbf{0}\}.$$

Assume that

$$\nu_1 \|\mathbf{y}\|^2 \leq (\mathcal{A}\mathbf{y}, \mathbf{y}) \leq \nu_2 \|\mathbf{y}\|^2, \quad \mathbf{y} \in \mathbf{U},$$

Let the operator \mathbf{B} possesses the following property: there exists a constant α such that for any

$$\mathbf{g} \in \text{Im } \mathbf{B} := \{\mathbf{z} \in \mathbf{H} \mid \exists \mathbf{v} \in \mathbf{W}_0 : \mathbf{B}\mathbf{v} = \mathbf{z}\}$$

one can find $\mathbf{u}_g \in \mathbf{W}_0$ such that

$$\mathbf{B}\mathbf{u}_g = \mathbf{g} \quad \text{and} \quad \|\mathbf{u}_g\|_{\mathbf{W}} \leq \alpha \|\mathbf{g}\|.$$

Note that such a condition is a generalization of Lemma 2.

Under the above assumption we obtain an estimate of the deviation from \mathbf{u} .

Estimate of the deviation from \mathbf{u}

$$\begin{aligned} \|\boldsymbol{\Lambda}(\mathbf{u} - \widehat{\mathbf{v}})\| &\leq \\ &\leq 2\sqrt{\nu_2}\alpha\|\mathbf{B}\widehat{\mathbf{v}}\| + \|\mathcal{A}\boldsymbol{\Lambda}\widehat{\mathbf{v}} - \mathbf{y}\|_* + \frac{1}{\sqrt{\nu_1}}\|\mathbf{f} + \boldsymbol{\Lambda}^*\mathbf{y} - \mathbf{B}^*\mathbf{q}\|. \end{aligned}$$

where $\|\mathbf{y}\| := (\mathcal{A}\mathbf{y}, \mathbf{y})^{1/2}$, $\|\mathbf{y}\|_* := (\mathcal{A}^{-1}\mathbf{y}, \mathbf{y})^{1/2}$. We see that the terms of the estimate present errors in the basic relations

$$\begin{cases} \langle \boldsymbol{\Lambda}^*\boldsymbol{\sigma} + \mathbf{f} - \mathbf{B}^*\mathbf{p}, \mathbf{w} \rangle = 0 & \forall \mathbf{w} \in \mathbf{V}_0, \\ \boldsymbol{\sigma} = \mathcal{A}\boldsymbol{\Lambda}\mathbf{u}, \\ \mathbf{B}\mathbf{v} = 0. \end{cases}$$

For the Stokes problem $\mathbf{A}\mathbf{v} = \nabla\mathbf{v}$, $\mathcal{A} = \nu\mathbb{I}$, where \mathbb{I} denotes the identity operator and $\mathbf{B}\mathbf{v} = -\mathbf{div}\mathbf{v}$. It is easy to see that in this case $\nu_1 = \nu_2 = \nu$,

$$\|\mathcal{A}\mathbf{A}\hat{\mathbf{v}} - \mathbf{y}\|_* = \frac{1}{\sqrt{\nu}} \|\nu\nabla\mathbf{v} - \mathbf{y}\|.$$

Since $\|\mathbf{A}(\mathbf{u} - \hat{\mathbf{v}})\| = \sqrt{\nu}\|\mathbf{A}(\mathbf{u} - \hat{\mathbf{v}})\|$, we find that the general estimate coincides with (262).

NONLINEAR MODELS IN THE THEORY OF FLUIDS

Bingham fluids

In these models

$$\mathbf{W}(\varepsilon) = \mu \left| \varepsilon^{\mathbf{D}} \right|^2 + \sqrt{2} \mathbf{K}_* \left| \varepsilon^{\mathbf{D}} \right| = \mathbf{W}_1(\varepsilon) + \mathbf{W}_2(\varepsilon) \quad (265)$$

Here \mathbf{W}_1 is the Newtonian potential of a **viscous fluid** and \mathbf{W}_2 is the **plastic** potential. By the Moreau–Rockafellar theorem we have

$$\partial \mathbf{W}(\varepsilon) = \partial \mathbf{W}_1(\varepsilon) + \partial \mathbf{W}_2(\varepsilon). \quad (266)$$

If $\left| \varepsilon^{\mathbf{D}} \right| \neq 0$ then the relation reads

$$\partial \mathbf{W}(\varepsilon) = 2\mu \varepsilon^{\mathbf{D}} + \frac{\sqrt{2} \mathbf{K}_*}{\left| \varepsilon^{\mathbf{D}} \right|} \varepsilon^{\mathbf{D}}. \quad (267)$$

Let us define $\partial \mathbf{W}(\varepsilon)$ for the case $|\varepsilon^{\mathbf{D}}| = \mathbf{0}$. Evidently, $\mathbf{W}_1(\varepsilon) = \mathbf{0}$. To find $\partial \mathbf{W}_2$ we recall that by definition, $\tau \in \partial W_2(\varepsilon)$ iff

$$\sqrt{2}\mathbf{K}_* \left| \mathfrak{a}^{\mathbf{D}} \right| - \sqrt{2}\mathbf{K}_* \left| \varepsilon^{\mathbf{D}} \right| \geq \tau : (\mathfrak{a} - \varepsilon).$$

Note that

$$\tau : (\mathfrak{a} - \varepsilon) = \frac{1}{n} \mathbf{tr} \tau (\mathbf{tr} \mathfrak{a} - \mathbf{tr} \varepsilon) + \tau^{\mathbf{D}} : \mathfrak{a}^{\mathbf{D}},$$

Therefore,

$$\sqrt{2}\mathbf{K}_* \left| \mathfrak{a}^{\mathbf{D}} \right| \geq \frac{1}{n} \mathbf{tr} \tau (\mathbf{tr} \mathfrak{a} - \mathbf{tr} \varepsilon) + \tau^{\mathbf{D}} : \mathfrak{a}^{\mathbf{D}} \quad \forall \mathfrak{a}^{\mathbf{D}}.$$

Let $|\boldsymbol{\varepsilon}^D| = \mathbf{0}$, then

$$\mathbf{0} \geq \frac{1}{n} \operatorname{tr} \boldsymbol{\tau} (\operatorname{tr} \boldsymbol{\varepsilon} - \operatorname{tr} \boldsymbol{\varepsilon}).$$

From here, it follows that

$$\operatorname{tr} \boldsymbol{\tau} = \mathbf{0}.$$

Therefore,

$$\sqrt{2\mathbf{K}_*} |\boldsymbol{\varepsilon}^D| \geq \boldsymbol{\tau}^D : \boldsymbol{\varepsilon}^D \quad \forall \boldsymbol{\varepsilon}^D.$$

Without a loss of generality we may assume that $|\boldsymbol{\varepsilon}^D| = \mathbf{1}$. To maximize the right-hand side we take $\boldsymbol{\varepsilon}^D := \frac{\boldsymbol{\tau}^D}{|\boldsymbol{\tau}^D|}$. Then, we observe that

$$\sqrt{2\mathbf{K}_*} \geq |\boldsymbol{\tau}^D|.$$

Hence, $\partial \mathbf{W}_2(\varepsilon)$ for $|\varepsilon^D| = \mathbf{0}$ consists of all τ such that $\tau = \tau^D$ and

$$|\tau^D| \leq \sqrt{2}K_*. \quad (268)$$

This relation reflects the physical nature of a **viscoplastic fluid**. In the stagnation zone the deviatoric part of its stress is not uniquely defined and can be any provides (268) is satisfied. Later we will see the examples exposing such a behavior of a solution.

Bingham flow in a basin

We consider a flow problem in a fixed domain, i.e., $\Omega_t = \Omega$ and the "no-slip" conditions $\mathbf{u} = \mathbf{u}_0$ are imposed on the boundary.

Classical solution is defined as $(\mathbf{u}, \boldsymbol{\sigma}, \mathbf{p})$ such that

$$\mathbf{u} \in \mathbf{C}(\bar{\Omega}; \mathbb{R}^d) \cap \mathbf{W}_2^1(\Omega; \mathbb{R}^d), \quad (269)$$

$$\boldsymbol{\sigma} \in \mathbf{C}(\bar{\Omega}; \mathbf{M}_s^{n \times n}) \cap \mathbf{W}_2^1(\Omega; \mathbf{M}_s^{n \times n}), \quad (270)$$

$$\mathbf{p} \in \mathbf{C}(\Omega), \quad (271)$$

$$\frac{\partial \mathbf{u}}{\partial \mathbf{t}} \in \mathbf{L}_2(\Omega; \mathbb{R}^d), \quad (272)$$

$$\mathbf{f} \in \mathbf{L}_2(\Omega; \mathbb{R}^d), \quad (273)$$

Consider Dirichlét boundary conditions.

Let $\mathbf{v} \in \overset{\circ}{\mathbf{H}}^1(\Omega; \mathbb{R}^d) = \{ \mathbf{v} \in \mathbf{H}^1(\Omega; \mathbb{R}^d) \mid \mathbf{v} = \mathbf{0} \text{ on } \partial\Omega \}$. Multiply

$$\frac{\partial \mathbf{u}}{\partial t} + \mathbf{u}_i \frac{\partial \mathbf{u}}{\partial \mathbf{x}_i} - \mathbf{div} \boldsymbol{\sigma} = \mathbf{f}$$

by \mathbf{v} and integrate by parts.

$$\int_{\Omega} \left(\frac{\partial \mathbf{u}}{\partial t} \cdot \mathbf{v} + \mathbf{u}_i \frac{\partial \mathbf{u}}{\partial \mathbf{x}_i} \cdot \mathbf{v} - \mathbf{v} \mathbf{div} \boldsymbol{\sigma} \right) \mathbf{d}\mathbf{x} = \int_{\Omega} \mathbf{f} \mathbf{v} \mathbf{d}\mathbf{x},$$

$$\int_{\Omega} \left(\frac{\partial \mathbf{u}}{\partial t} + \mathbf{u}_i \frac{\partial \mathbf{u}}{\partial \mathbf{x}_i} \right) \cdot \mathbf{v} \mathbf{d}\mathbf{x} + \int_{\Omega} \boldsymbol{\sigma} : \boldsymbol{\varepsilon}(\mathbf{v}) \mathbf{d}\mathbf{x} = \int_{\Omega} \mathbf{f} \mathbf{v} \mathbf{d}\mathbf{x}.$$

This relation may be rewritten as

$$\int_{\Omega} \left(\frac{\partial \mathbf{u}}{\partial t} + \mathbf{u}_i \frac{\partial \mathbf{u}}{\partial \mathbf{x}_i} \right) \cdot \mathbf{v} \mathbf{d}\mathbf{x} + \int_{\Omega} (\boldsymbol{\sigma} + \mathbf{p}\mathbf{l}) : \boldsymbol{\varepsilon}(\mathbf{v}) \mathbf{d}\mathbf{x} = \int_{\Omega} (\mathbf{f} \mathbf{v} + \mathbf{p} \mathbf{div} \mathbf{v}) \mathbf{d}\mathbf{x} \quad \forall \mathbf{v} \in \overset{\circ}{\mathbf{H}}^1(\Omega; \mathbb{R}^d) \quad (274)$$

By the constitutive relations,

$$\boldsymbol{\sigma} + \mathbf{p}\mathbf{l} = \partial\mathbf{W}_1 + \partial\mathbf{W}_2 = 2\mu\boldsymbol{\varepsilon}^D(\mathbf{u}) + \partial\mathbf{W}_2(\boldsymbol{\varepsilon}(\mathbf{u}))$$

what means that

$$\boldsymbol{\sigma} + \mathbf{p}\mathbf{l} - 2\mu\boldsymbol{\varepsilon}^D(\mathbf{u}) \in \partial\mathbf{W}_2(\boldsymbol{\varepsilon}(\mathbf{u}))$$

The latter means that

$$\begin{aligned} \sqrt{2}\mathbf{k}_* \left| \boldsymbol{\varepsilon}^D \right| - \sqrt{2}\mathbf{k}_* \left| \boldsymbol{\varepsilon}^D(\mathbf{u}) \right| &\geq \\ &\geq (\boldsymbol{\sigma} + \mathbf{p}\mathbf{l} - 2\mu\boldsymbol{\varepsilon}^D(\mathbf{u})) : (\boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}(\mathbf{u})) \quad \forall \boldsymbol{\varepsilon} \in \mathbf{M}_s^{n \times n} \end{aligned}$$

Set $\mathfrak{a} = \varepsilon(\mathbf{u} + \mathbf{v})$. Then,

$$\sqrt{2\mathbf{k}_*} \left| \mathfrak{a}^D \right| - \sqrt{2\mathbf{k}_*} \left| \varepsilon^D(\mathbf{u}) \right| \geq (\boldsymbol{\sigma} + \mathbf{p}\mathbf{l} - 2\mu\varepsilon^D(\mathbf{u})) : \varepsilon(\mathbf{v}).$$

and we obtain

$$\begin{aligned} \int_{\Omega} (\boldsymbol{\sigma} + \mathbf{p}\mathbf{l}) : \varepsilon(\mathbf{v}) \, dx &\leq \\ &\leq \int_{\Omega} (2\mu\varepsilon^D(\mathbf{u}) : \varepsilon^D(\mathbf{v}) + \sqrt{2\mathbf{k}_*} \left| \mathfrak{a}^D \right| - \sqrt{2\mathbf{k}_*} \left| \varepsilon^D(\mathbf{u}) \right|) \, dx \end{aligned}$$

Now (274) yields inequality

$$\begin{aligned} \int_{\Omega} \left(\frac{\partial \mathbf{u}}{\partial \mathbf{t}} + \mathbf{u}_i \frac{\partial \mathbf{u}}{\partial \mathbf{x}_i} \right) \cdot \mathbf{v} \, dx + \int_{\Omega} (2\mu\varepsilon^D(\mathbf{u}) : \varepsilon^D(\mathbf{v}) + \\ + \sqrt{2\mathbf{k}_*} (\left| \mathfrak{a}^D \right| - \left| \varepsilon^D(\mathbf{u}) \right|)) \, dx \geq \int_{\Omega} (\mathbf{f}\mathbf{v} + \mathbf{p} \operatorname{div} \mathbf{v}) \, dx \quad \forall \mathbf{v} \in \mathbf{H}^1(\Omega; \mathbb{R}^d) \end{aligned}$$

If a trial function is taken from a narrower set

$$\mathbf{v} = \mathbf{w} + \mathbf{u}_0 - \mathbf{u} \quad , \quad \mathbf{w}, \mathbf{u}_0 \in \mathring{\mathbf{J}}_2^1(\Omega)$$

so that

$$\mathbf{div} \mathbf{v} = \mathbf{0}, \quad \mathbf{\varepsilon} = \varepsilon(\mathbf{w} + \mathbf{u}_0)$$

Then (275) comes in a simpler form

$$\begin{aligned} & \int_{\Omega} \left\{ \left(\frac{\partial \mathbf{u}}{\partial t} + \mathbf{u}_i \frac{\partial \mathbf{u}}{\partial x_i} \right) \cdot (\mathbf{w} + \mathbf{u}_0 - \mathbf{u}) + 2\mu \varepsilon^D(\mathbf{u}) : \varepsilon^D(\mathbf{w} + \mathbf{u}_0 - \mathbf{u}) + \right. \\ & \left. + \sqrt{2} \mathbf{k}_* \left(\left| \varepsilon^D(\mathbf{w} + \mathbf{u}_0) \right| - \left| \varepsilon^D(\mathbf{u}) \right| \right) \right\} \mathbf{d}\mathbf{x} \geq \int_{\Omega} \mathbf{f} \cdot (\mathbf{w} + \mathbf{u}_0 - \mathbf{u}) \mathbf{d}\mathbf{x} \quad \forall \mathbf{w} \in \mathring{\mathbf{J}}_2^1(\Omega). \end{aligned} \quad (275)$$

Denote $\tilde{\mathbf{u}} = \mathbf{u}_0 + \mathbf{w}$, then we arrive at the variational inequality that describes the motion of an elastoplastic media

$$\begin{aligned} \int_{\Omega} \left\{ \left(\frac{\partial \mathbf{u}}{\partial t} + \mathbf{u}_i \frac{\partial \mathbf{u}}{\partial x_i} \right) \cdot (\tilde{\mathbf{u}} - \mathbf{u}) + \right. \\ \left. + 2\mu \varepsilon^D(\mathbf{u}) : \varepsilon^D(\tilde{\mathbf{u}} - \mathbf{u}) + \sqrt{2} \mathbf{k}_* \left(\left| \varepsilon^D(\tilde{\mathbf{u}}) \right| - \left| \varepsilon^D(\mathbf{u}) \right| \right) \right\} \mathbf{d}\mathbf{x} \geq \\ \geq \int_{\Omega} \mathbf{f} \cdot (\tilde{\mathbf{u}} - \mathbf{u}) \mathbf{d}\mathbf{x} \quad \forall \tilde{\mathbf{u}} \in \mathbf{u}_0 + \mathring{\mathbf{J}}_2^1(\Omega) \quad (276) \end{aligned}$$

Stationary flow

Stationary model is a particular case of the above that arises if $\partial \mathbf{u} / \partial \mathbf{t} = \mathbf{0}$.

Here, we need to find $\mathbf{u} \in \mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0$ such that

$$\int_{\Omega} \left\{ \mathbf{u}_i \frac{\partial \mathbf{u}}{\partial x_i} \cdot (\mathbf{v} - \mathbf{u}) + 2\mu \varepsilon(\mathbf{u}) : \varepsilon(\mathbf{v} - \mathbf{u}) + \sqrt{2} \mathbf{k}_* (|\varepsilon(\mathbf{v})| - |\varepsilon(\mathbf{u})|) \right\} \mathbf{d}\mathbf{x} \geq \int_{\Omega} \mathbf{f} \cdot (\mathbf{v} - \mathbf{u}) \mathbf{d}\mathbf{x} \quad \forall \mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0. \quad (277)$$

If, in addition, \mathbf{v} is small (slow flows) then we arrive at the problem

$$\begin{aligned} \int_{\Omega} \left\{ 2\mu \varepsilon(\mathbf{u}) : \varepsilon(\mathbf{v} - \mathbf{u}) + \sqrt{2} \mathbf{k}_* (|\varepsilon(\mathbf{v})| - |\varepsilon(\mathbf{u})|) \right\} \mathbf{d}\mathbf{x} &\geq & (278) \\ &\geq \int_{\Omega} \mathbf{f} \cdot (\mathbf{v} - \mathbf{u}) \mathbf{d}\mathbf{x} \quad \forall \mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0 \end{aligned}$$

Existence and uniqueness

Consider the problem: find $\mathbf{u} \in \mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0$ such that

$$\begin{aligned} \int_{\Omega} 2\mu \boldsymbol{\varepsilon}(\mathbf{u}) \cdot \boldsymbol{\varepsilon}(\mathbf{v} - \mathbf{u}) + \sqrt{2\mathbf{K}_*} (|\boldsymbol{\varepsilon}(\mathbf{v})| - |\boldsymbol{\varepsilon}(\mathbf{u})|) \mathbf{d}\mathbf{x} &\geq \\ &\geq \int_{\Omega} \mathbf{f} \cdot (\mathbf{v} - \mathbf{u}) \mathbf{d}\mathbf{x} \quad \forall \mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0, \quad (279) \end{aligned}$$

where $\mathbf{f} \in \mathbf{L}_2(\Omega, \mathbb{R}^d)$, $\mathbf{u}_0 \in \mathring{\mathbf{J}}_2^1(\Omega)$.

Uniqueness of the solution is easy to prove

Regularity of weak solutions

Regularity of weak solutions in the theory of Non Newtonian viscous fluids was investigated in the works of M. Fuchs, G. Seregin, M. Bildhauer and other authors. Readers can find a consequent exposition of these results in the book

M. Fuchs and G.A. Seregin. *Variational methods for problems from plasticity theory and for generalized Newtonian fluids.* Lect. Notes in Mathematics 1749, Springer-Verlag, Berlin (2000)

Let $\mathbf{u}_1, \mathbf{u}_2 \in \mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0$ be two different solutions. Then

$$\begin{aligned} & \int_{\Omega} \{ \mu \boldsymbol{\varepsilon}(\mathbf{u}_1) \cdot \boldsymbol{\varepsilon}(\mathbf{v} - \mathbf{u}_1) + \\ & \quad + \sqrt{2} \mathbf{K}_* (|\boldsymbol{\varepsilon}(\mathbf{v})| - |\boldsymbol{\varepsilon}(\mathbf{u}_1)|) - \mathbf{f} \cdot (\mathbf{v} - \mathbf{u}_1) \} \mathbf{d}\mathbf{x} \geq \mathbf{0} \\ & \int_{\Omega} \{ \mu \boldsymbol{\varepsilon}(\mathbf{u}_2) \cdot \boldsymbol{\varepsilon}(\mathbf{v} - \mathbf{u}_2) + \\ & \quad + \sqrt{2} \mathbf{K}_* (|\boldsymbol{\varepsilon}(\mathbf{v})| - |\boldsymbol{\varepsilon}(\mathbf{u}_2)|) - \mathbf{f} \cdot (\mathbf{v} - \mathbf{u}_2) \} \mathbf{d}\mathbf{x} \geq \mathbf{0} \end{aligned}$$

Set $\mathbf{v} = \mathbf{u}_2$ in the first and $\mathbf{v} = \mathbf{u}_1$ in the second.

$$\int_{\Omega} \mu \boldsymbol{\varepsilon}(\mathbf{u}_1 - \mathbf{u}_2) \cdot \boldsymbol{\varepsilon}(\mathbf{u}_2 - \mathbf{u}_1) \mathbf{d}\mathbf{x} \geq \mathbf{0}$$

Thus,

$$\int_{\Omega} |\boldsymbol{\varepsilon}(\mathbf{u}_1 - \mathbf{u}_2)|^2 \mathbf{d}\mathbf{x} \leq \mathbf{0}$$

Existence

Existence follows from that the problem is equivalent to the variational problem

$$I(\mathbf{u}) = \inf \{I(\mathbf{v}) : \mathbf{v} \in \mathring{J}_2^1(\Omega) + \mathbf{u}_0\}, \quad (280)$$

where

$$I(\mathbf{v}) = \int_{\Omega} \{ \mu |\varepsilon(\mathbf{u})|^2 + \sqrt{2}K_* |\varepsilon(\mathbf{v})| - \mathbf{f} \cdot \mathbf{v} \} \mathbf{d}\mathbf{x}$$

This equivalence follows from the well known result (see, e.g., [J.-L. Lions and G. Duvaut. Inequalities in mechanics and physics](#)).

Theorem

Variational problem

$$\inf_{\mathbf{v} \in \mathbf{K}} \mathbf{J}(\mathbf{v}) = \mathbf{J}(\mathbf{v}_*); \quad \mathbf{J}(\mathbf{v}) = \mathbf{J}_0(\mathbf{v}) + \mathbf{j}(\mathbf{v}), \quad (281)$$

where \mathbf{K} is a convex closed set, $\mathbf{J}_0 : \mathbf{V} \rightarrow \mathbb{R}$ is a convex differentiable functional, $\mathbf{j} : \mathbf{V} \rightarrow \mathbb{R}$ is a convex functional, \mathbf{V} has the same solution as the variational inequality

$$(\mathbf{J}'_0(\mathbf{v}_*), \mathbf{v}_* - \mathbf{v}) + \mathbf{j}(\mathbf{v}_*) - \mathbf{j}(\mathbf{v}) \leq 0, \quad \forall \mathbf{v} \in \mathbf{K} \quad (282)$$

In our case

$$\mathbf{J}_0(\mathbf{v}) = \int_{\Omega} \mu |\varepsilon(\mathbf{u})|^2 \, d\mathbf{x} - \int_{\Omega} \mathbf{f} \cdot \mathbf{v} \, d\mathbf{x}$$

$$\mathbf{j}(\mathbf{v}) = \sqrt{2\mathbf{K}_*} \int_{\Omega} |\varepsilon(\mathbf{u})| \, d\mathbf{x}$$

$$\mathbf{K} = \overset{\circ}{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0, \quad \mathbf{V} = \mathbf{H}^1(\Omega, \mathbb{R}^n)$$

and we arrive at

$$\mathbf{l}(\mathbf{v}) = \int_{\Omega} \{ \mu |\varepsilon(\mathbf{u})|^2 + \sqrt{2\mathbf{K}_*} |\varepsilon(\mathbf{v})| - \mathbf{f} \cdot \mathbf{v} \} \, d\mathbf{x}$$

which is **convex, continuous, and coercive**. These properties imply existence.

Approximate solutions

Let $\mathbf{V}_h \subset \mathring{\mathbf{J}}_2^1(\Omega)$ and $\dim \mathbf{V}_h < +\infty$, find $\mathbf{u}_h \in \mathbf{V}_h + \mathbf{u}_0$,

$$\int_{\Omega} \mu \varepsilon(\mathbf{u}_h) \cdot (\varepsilon(\mathbf{v}_h) - \varepsilon(\mathbf{u}_h)) \, \mathbf{d}\mathbf{x} + \int_{\Omega} \sqrt{2} \mathbf{K}_* (|\varepsilon(\mathbf{v}_h)| - |\varepsilon(\mathbf{u}_h)|) \, \mathbf{d}\mathbf{x} \geq \int_{\Omega} \mathbf{f} \cdot (\mathbf{v}_h - \mathbf{u}_h) \, \mathbf{d}\mathbf{x} \quad \forall \mathbf{v}_h \in \mathbf{V}_h + \mathbf{u}_0 \quad (283)$$

Variational formulation: find $\mathbf{u}_h \in \mathbf{V}_h + \mathbf{u}_0$ such that

$$\mathbf{l}(\mathbf{u}_h) = \inf \{ \mathbf{l}(\mathbf{v}_h) : \mathbf{v}_h \in \mathbf{V}_h + \mathbf{u}_0 \} \quad (284)$$

Serious difficulty is the condition

$$\mathbf{div} \mathbf{v}_h = \mathbf{0}.$$

Example. Antiplane flow in a pipe

Consider a long tube of the cross-section Ω with the income and outcome pressure:

$$\mathbf{p}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{0}, t) = \mathbf{0}, \quad \mathbf{p}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{L}, t) = -c\mathbf{L} \quad (285)$$

On the side surface

$$\mathbf{u} \big|_{\partial\Omega \times [0, L]} = \mathbf{0}, \quad (286)$$

We are looking for a solution of the form

$$\mathbf{u}_1 = \mathbf{0}, \quad \mathbf{u}_2 = \mathbf{0}, \quad \mathbf{u}_3 = \mathbf{w}(\mathbf{x}_1, \mathbf{x}_2, t)$$

$$\varepsilon(\mathbf{u}) = \begin{pmatrix} 0 & 0 & \frac{1}{2}\mathbf{w}_{,1} \\ 0 & 0 & \frac{1}{2}\mathbf{w}_{,2} \\ \frac{1}{2}\mathbf{w}_{,1} & \frac{1}{2}\mathbf{w}_{,2} & 0 \end{pmatrix}$$

Since

$$\varepsilon(\mathbf{u}) = \begin{pmatrix} 0 & 0 & \frac{1}{2}\mathbf{w},1 \\ 0 & 0 & \frac{1}{2}\mathbf{w},2 \\ \frac{1}{2}\mathbf{w},1 & \frac{1}{2}\mathbf{w},2 & 0 \end{pmatrix}$$

we have

$$|\varepsilon(\mathbf{u})|^2 = \frac{1}{2} |\nabla \mathbf{w}|^2, \quad \frac{\partial \mathbf{u}}{\partial \mathbf{x}_1} = (\mathbf{0}, \mathbf{0}, \mathbf{w},1),$$

$$\frac{\partial \mathbf{u}}{\partial \mathbf{x}_2} = (\mathbf{0}, \mathbf{0}, \mathbf{w},2), \quad \frac{\partial \mathbf{u}}{\partial \mathbf{x}_3} = (\mathbf{0}, \mathbf{0}, \mathbf{0}), \quad \mathbf{u}_i \frac{\partial \mathbf{u}}{\partial \mathbf{x}_i} = \mathbf{0}$$

$$\mathbf{u}_i \frac{\partial \mathbf{u}}{\partial \mathbf{x}_i} = \mathbf{0}, \quad \operatorname{div} \mathbf{u} = \mathbf{0} \quad (287)$$

Remark

If set in (277) $\mathbf{v} = (\mathbf{0}, \mathbf{0}, \tilde{\mathbf{w}})$, where $\tilde{\mathbf{w}} \in \overset{\circ}{\mathbf{H}}^1(\Omega)$, $\tilde{\mathbf{w}} = \tilde{\mathbf{w}}(x_1, x_2)$, $\mathbf{u}_0 = \mathbf{0}$ and $\mathbf{f} = (\mathbf{0}, \mathbf{0}, \mathbf{c})$, then we arrive at the problem

$$\begin{aligned} \int_{\Omega} \{ \mu \nabla \mathbf{w} \cdot (\nabla \tilde{\mathbf{w}} - \nabla \mathbf{w}) + \mathbf{k}_*(|\nabla \tilde{\mathbf{w}}| - |\nabla \mathbf{w}|) \} \mathbf{d}\mathbf{x} &\geq \\ &\geq \int_{\Omega} \mathbf{c}(\tilde{\mathbf{w}} - \mathbf{w}) \mathbf{d}\mathbf{x} \quad \forall \tilde{\mathbf{w}} \in \overset{\circ}{\mathbf{H}}^1(\Omega). \end{aligned}$$

We see that stationary slow flow of a viscoplastic fluid in a pipe is described by a **variational inequality**.

In this case,

$$\boldsymbol{\sigma} + \rho \mathbf{l} = \boldsymbol{\tau} = \begin{pmatrix} 0 & 0 & \tau_{31}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{t}) \\ 0 & 0 & \tau_{32}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{t}) \\ \tau_{31}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{t}) & \tau_{32}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{t}) & 0 \end{pmatrix}$$

and $\boldsymbol{\tau} = \boldsymbol{\tau}^D$.

Let $\rho_0 = 1$. Take the third equation of the motion:

$$\frac{\partial \mathbf{w}}{\partial \mathbf{t}} = \sigma_{31,1} + \sigma_{32,2} + \sigma_{33,3}, \quad (288)$$

Here $\boldsymbol{\sigma} = \boldsymbol{\tau} - \rho \mathbf{l}$ and $\sigma_{33} = -\rho$.

Therefore, we have

$$\frac{\partial \mathbf{w}}{\partial \mathbf{t}} = \boldsymbol{\sigma}_{31,1} + \boldsymbol{\sigma}_{32,2} - \frac{\partial \mathbf{p}}{\partial \mathbf{x}_3} \quad (289)$$

Since $\boldsymbol{\sigma}_{31}$, $\boldsymbol{\sigma}_{32}$, and \mathbf{w} depend only on x_1, x_2 , we rewrite 289 as follows:

$$\frac{\partial \mathbf{p}}{\partial \mathbf{x}_3} = \mathbf{a}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{t})$$

Thus,

$$\mathbf{p}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{t}) = \mathbf{a}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{t}) + \mathbf{x}_3 \mathbf{b}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{t}).$$

Boundary conditions (285) say that

$$\begin{aligned} \mathbf{p}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{0}, \mathbf{t}) &= \mathbf{a}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{t}) = \mathbf{0}, \\ \mathbf{p}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{L}, \mathbf{t}) &= \mathbf{Lb}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{t}) = -\mathbf{cL} \implies \mathbf{b} = -\mathbf{c}. \end{aligned}$$

Hence, the pressure is found as $\mathbf{p} = -\mathbf{Cx}_3$.

Now the motion is governed by the single equation

$$\frac{\partial \mathbf{w}}{\partial \mathbf{t}} = \frac{\partial \sigma_{31}}{\partial \mathbf{x}_1} + \frac{\partial \sigma_{32}}{\partial \mathbf{x}_2} + \mathbf{c} \quad (290)$$

Recall that

$$|\varepsilon(\mathbf{u})|^2 = \frac{1}{2} |\nabla \mathbf{w}|^2$$

Rewrite the constitutive relation

$$\sigma + \mathbf{pl} = \begin{cases} 2\mu\varepsilon(\mathbf{u}) + \sqrt{2}\mathbf{K}_* \frac{\varepsilon(\mathbf{u})}{|\varepsilon(\mathbf{u})|} & \text{for } |\varepsilon(\mathbf{u})| > 0 \\ \tau = \tau^D, |\tau^D| \leq \sqrt{2}\mathbf{K}_* & \text{for } |\varepsilon(\mathbf{u})| = 0 \end{cases}$$

in terms of the problem considered.

We have

$$\sigma_{3i} = \begin{cases} (2\mu + \sqrt{2}K_* \frac{1}{\frac{1}{\sqrt{2}} |\nabla \mathbf{w}|}) \frac{1}{2} \mathbf{w}_{,i} & |\nabla \mathbf{w}| > 0, \\ \tau_{3i}, & |\tau^D| \leq \sqrt{2}K_* \quad |\nabla \mathbf{w}| = 0, \end{cases} \quad i = 1, 2. \quad (291)$$

where $|\tau^D|^2 = 2\tau_{31}^2 + 2\tau_{32}^2$. Therefore,

$$\sigma_{3i} = \begin{cases} (\mu + \frac{K_*}{|\nabla \mathbf{w}|}) \mathbf{w}_{,i}, & \text{for } |\nabla \mathbf{w}| > 0, \\ \tau_{3i}, \quad \sqrt{\tau_{31}^2 + \tau_{32}^2} \leq K_*, & \text{for } |\nabla \mathbf{w}| = 0. \end{cases} \quad (292)$$

We observe that

$$|\boldsymbol{\sigma}_3| := \sqrt{(\sigma_{31}^2 + \sigma_{32}^2)} = \begin{cases} \mu |\nabla \mathbf{w}| + \mathbf{K}_*, & \text{for } |\nabla \mathbf{w}| > \mathbf{0}, \\ \sqrt{\tau_{31}^2 + \tau_{32}^2} \leq \mathbf{K}_*, & \text{for } |\nabla \mathbf{w}| = \mathbf{0} \end{cases} .$$

Now (292), the equation

$$\frac{\partial \mathbf{w}}{\partial \mathbf{t}} = \mathbf{div} \boldsymbol{\sigma}_3 + \mathbf{c} .$$

and the condition

$$\mathbf{w} |_{\partial \Omega} = \mathbf{0} \tag{293}$$

describes the solution.

Cylindrical pipe

Consider stationary flow in a cylindrical pipe.

$$\boldsymbol{\sigma}_3 := \boldsymbol{\eta} = (\eta_\rho, \eta_\varphi), \quad \eta_\rho = \boldsymbol{\sigma}_3 \boldsymbol{e}_\rho, \quad \eta_\varphi = \boldsymbol{\sigma}_3 \boldsymbol{e}_\varphi.$$

Problem is axisymmetric and

$$\boldsymbol{\sigma}_3 \boldsymbol{e}_\varphi = 0, \quad \frac{\partial \boldsymbol{w}}{\partial \varphi} = 0.$$

Problem is stationary:

$$\frac{\partial \boldsymbol{w}}{\partial t} = 0.$$

In the axisymmetric case

$$\mathbf{div}(\boldsymbol{\sigma}_3) = \mathbf{div} \boldsymbol{\eta} = \frac{\partial \eta_\rho}{\partial \rho} + \frac{\eta_\rho}{\rho},$$

Now, the motion equation reads

$$\frac{\partial \eta_\rho}{\partial \rho} + \frac{\eta_\rho}{\rho} + \mathbf{c} = \mathbf{0} \implies \frac{\mathbf{1}}{\rho} \frac{\partial(\rho \eta_\rho)}{\partial \rho} + \mathbf{c} = \mathbf{0} \quad (294)$$

$$\eta_\rho = \begin{cases} \left(\mu + \frac{\mathbf{K}_*}{\left| \frac{\partial \mathbf{w}}{\partial \rho} \right|} \right) \frac{\partial \mathbf{w}}{\partial \rho} & \text{for } \left| \frac{\partial \mathbf{w}}{\partial \rho} \right| > \mathbf{0}, \\ \tilde{\eta}_\rho, \left| \tilde{\eta}_\rho \right| \leq \mathbf{K}_* & \text{for } \left| \frac{\partial \mathbf{w}}{\partial \rho} \right| = 0 \end{cases} \quad (295)$$

From the viewpoint of physics it would be natural to await the following behavior of the media:

1. If the "pressure grade" c is small, then there is no motion at all: plastic properties of the media dominate and does not allow any motion.
2. If c becomes large enough, then the motion starts in places where the effective stresses achieve the critical value, i.e. near the boundary. Central part moves as a solid body.

Formal analysis confirms these expectations. Integrate (294), then we observe that

$$\eta_\rho = -\frac{c}{2}\rho + \frac{C_2}{\rho} \quad (296)$$

$\sigma_{3\rho}$ at the center must be finite, therefore $C_2 = 0$ and (296) reads

$$\eta_\rho = -\frac{c}{2}\rho \quad (297)$$

For which c no motion arise?

$$|\eta_\rho(\rho)| \leq |\eta_\rho(\mathbf{R})| = \frac{c}{2}\mathbf{R} \leq \mathbf{K}_*,$$

Hence, these values are: $c \leq \frac{2\mathbf{K}_*}{\mathbf{R}}$.

Indeed, in such a case we have a function η_ρ that satisfies the equation and constitutive relations related to the branch $|\frac{\partial \mathbf{w}}{\partial \rho}| = 0$. Since $\mathbf{w}|_{\partial\Omega} = 0$ we conclude that $\mathbf{w} \equiv \mathbf{0}$.

Let $\mathbf{c} > \frac{2\mathbf{K}_*}{\mathbf{R}}$. We are looking for a solution such that

1. for $R_* \leq \rho \leq R$ the media is deformed and $|\frac{\partial \mathbf{w}}{\partial \rho}| > 0$
2. for $0 \leq \rho \leq R_*$ the media is rigid and $|\frac{\partial \mathbf{w}}{\partial \rho}| = 0$.

Then \mathbf{R}_* is defined by the relation

$$|\eta_\rho(\mathbf{R}_*)| = \frac{\mathbf{c}}{2}\mathbf{R}_* = \mathbf{K}_* \implies \mathbf{R}_* = \frac{2\mathbf{K}_*}{\mathbf{c}} < \mathbf{R}$$

For $\mathbf{R}_* \leq \rho \leq \mathbf{R}$ it should be $\frac{\partial \mathbf{w}}{\partial \rho} < 0$, so that

$$\eta_\rho = \left(\mu + \frac{\mathbf{K}_*}{\left| \frac{\partial \mathbf{w}}{\partial \rho} \right|} \right) \frac{d\mathbf{w}}{d\rho} = \mu \frac{d\mathbf{w}}{d\rho} - \mathbf{K}_* .$$

The latter relation leads to the conclusion that

$$\begin{aligned} \mu \frac{d\mathbf{w}}{d\rho} - \mathbf{K}_* &= -\frac{\mathbf{c}}{2}\rho \implies \mu \mathbf{w} - \mathbf{K}_* \rho = -\frac{\mathbf{c}}{4}\rho^2 + \mathbf{C}_3 \\ \mathbf{0} &= \mathbf{K}_* \mathbf{R} - \frac{\mathbf{c}}{4}\mathbf{R}^2 + \mathbf{C}_3 \implies \mathbf{C}_3 = \frac{\mathbf{c}}{4}\mathbf{R}^2 - \mathbf{K}_* \mathbf{R} . \end{aligned} \quad (298)$$

Consequently $\mu \mathbf{w} = \mathbf{K}_*(\rho - \mathbf{R}) - \frac{\mathbf{c}}{4}(\rho^2 - \mathbf{R}^2)$.

Therefore, we arrive at the conclusion that for $\mathbf{c} > \frac{2\mathbf{K}_*}{\mathbf{R}}$ the solution is as follows:

$$\mathbf{w} = \begin{cases} \frac{1}{\mu} \left[(\mathbf{R}_* - \mathbf{R}) - \frac{\mathbf{c}}{4} (\mathbf{R}_*^2 - \mathbf{R}^2) \right] & \mathbf{0} \leq \rho \leq \mathbf{R}_*, \\ \frac{1}{\mu} \left[(\rho - \mathbf{R}) - \frac{\mathbf{c}}{4} (\rho^2 - \mathbf{R}^2) \right] & \mathbf{R}_* \leq \rho \leq \mathbf{R}. \end{cases}$$

Thus,

$$\frac{d\mathbf{w}}{d\rho} = \mathbf{0} \quad \text{for } \mathbf{0} \leq \rho \leq \mathbf{R}_*$$

and

$$\frac{1}{\mu} \left[\mathbf{1} - \frac{\mathbf{c}}{2} \rho \right] \quad \text{for } \mathbf{R}_* \leq \rho \leq \mathbf{R}.$$

FUNCTIONAL A POSTERIORI ESTIMATES FOR GENERALIZED NEWTONIAN FLUIDS

For the considered class of problems it is convenient to derive estimates of deviations from exact solutions by the **variational** method. First, we establish the estimate

$$\frac{\nu}{2} \int_{\Omega} |\varepsilon(\mathbf{v} - \mathbf{u})|^2 \, d\mathbf{x} \leq \mathbf{J}(\mathbf{v}) - \mathbf{J}(\mathbf{u}), \quad (299)$$

where \mathbf{v} is an arbitrary function in $\mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0$. Indeed,

$$\begin{aligned} \mathbf{J}(\mathbf{v}) - \mathbf{J}(\mathbf{u}) = & \int_{\Omega} \left(\frac{\nu}{2} |\varepsilon(\mathbf{v} - \mathbf{u})|^2 + \nu \varepsilon(\mathbf{u}) : \varepsilon(\mathbf{v} - \mathbf{u}) + \right. \\ & \left. + \sqrt{2} \mathbf{K}_* (|\varepsilon(\mathbf{v})| - |\varepsilon(\mathbf{u})|) \right) d\mathbf{x} - \int_{\Omega} \mathbf{f} \cdot (\mathbf{v} - \mathbf{u}) \, d\mathbf{x}. \end{aligned}$$

To estimate $J(\mathbf{u})$ from below, we construct a set of variational problems whose functionals are defined on the functional class wider than $\mathring{J}_2^1(\Omega) + \mathbf{u}_0$. These problems we shall call "disturbed".

Let us define two functions $\mathbf{q} \in \mathring{\mathbf{L}}_2(\Omega)$ and $\boldsymbol{\tau}_2 \in \boldsymbol{\Sigma}$ and the functional

$$\bar{J}(\mathbf{v}) := \int_{\Omega} \left(\frac{\nu}{2} |\boldsymbol{\varepsilon}(\mathbf{v})|^2 + \boldsymbol{\tau}_2 : \boldsymbol{\varepsilon}(\mathbf{v}) - \psi^*(\boldsymbol{\tau}_2) - \mathbf{f} \cdot \mathbf{v} - \mathbf{q} \operatorname{div}(\mathbf{v} - \mathbf{u}_0) \right) \mathbf{d}\mathbf{x},$$

where $\psi^* : \mathbb{M}^{\mathbf{d} \times \mathbf{d}} \rightarrow \mathbb{R}$ is the functional conjugate to $\psi(\kappa) := \sqrt{2}K_*|\kappa|$ in the sense of Young–Fenchel, i.e.,

$$\psi^*(\kappa^*) = \sup_{\kappa \in \mathbb{M}^{\mathbf{d} \times \mathbf{d}}} \{ \kappa^* : \kappa - \psi(\kappa) \}.$$

Now, the following variational problem $\bar{\mathcal{P}}_{q,\tau_2}$ arises: Find $\bar{\mathbf{u}} \in \mathbf{V}_0 + \mathbf{u}_0$ such that

$$\bar{\mathbf{J}}(\bar{\mathbf{u}}) = \inf \bar{\mathcal{P}} := \inf_{\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0} \bar{\mathbf{J}}(\mathbf{v}).$$

Here, in principle, $\bar{\mathbf{u}} = \bar{\mathbf{u}}_{q,\tau_2}$. However, for the sake of simplicity we shall not do this assuming that the bar above means that a quantity depends on the above functions.

Problem $\bar{\mathcal{P}}$ is uniquely solvable and

$$\inf \bar{\mathcal{P}} \leq \inf \mathcal{P}. \quad (300)$$

Existence and uniqueness of $\bar{\mathbf{u}}$ follows from the properties of the convex functional $\bar{\mathbf{J}}$ and closed set $\mathbf{V}_0 + \mathbf{u}_0$. In accordance with the definition of ψ^* , for any $\mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0$ we have the inequality

$$\bar{\mathbf{J}}(\mathbf{v}) = \int_{\Omega} \left(\frac{\nu}{2} |\varepsilon(\mathbf{v})|^2 + \tau_2 : \varepsilon(\mathbf{v}) - \psi^*(\tau_2) \right) \mathbf{d}\mathbf{x} - \int_{\Omega} \mathbf{f} \cdot \mathbf{v} \mathbf{d}\mathbf{x} \leq \mathbf{J}(\mathbf{v}).$$

Therefore,

$$\inf_{\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0} \bar{\mathbf{J}}(\mathbf{v}) \leq \inf_{\mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0} \bar{\mathbf{J}}(\mathbf{v}) \leq \inf_{\mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0} \mathbf{J}(\mathbf{v}) = \inf \mathcal{P}.$$

Now, we observe that

$$\frac{\nu}{2} \int_{\Omega} |\varepsilon(\mathbf{v} - \mathbf{u})|^2 \, d\mathbf{x} \leq \mathbf{J}(\mathbf{v}) - \inf \bar{\mathcal{P}}. \quad (301)$$

However, the value of $\inf \bar{\mathcal{P}}$ is unknown! To overcome this difficulty we attract dual variational problem $\bar{\mathcal{P}}^*$.

If $\inf \bar{\mathcal{P}} = \sup \bar{\mathcal{P}}^*$, then $\inf \bar{\mathcal{P}}$ can be replaced by a lower estimate of $\sup \bar{\mathcal{P}}^*$. Estimates obtained on this way will depend on the functions $\boldsymbol{\tau}_2$ and \mathbf{q} and also on the variables of the dual problem. Note that there are different variational problem that may be viewed as dual to $\bar{\mathcal{P}}$. The problem is to find in this collection a proper variant that leads to estimates convenient for practice and having a good accuracy. For the class of problems considered, the following variant is possible.

Define the Lagrangian

$$\begin{aligned} \bar{\mathbf{L}}(\mathbf{v}; \boldsymbol{\tau}_1) := & \int_{\Omega} \left(\boldsymbol{\varepsilon}(\mathbf{v}) : (\boldsymbol{\tau}_1 + \boldsymbol{\tau}_2) - \frac{1}{2\nu} |\boldsymbol{\tau}_1|^2 - \psi^*(\boldsymbol{\tau}_2) \right) \mathbf{d}\mathbf{x} \\ & - \int_{\Omega} \mathbf{f} \cdot \mathbf{v} \mathbf{d}\mathbf{x} - \int_{\Omega} \mathbf{q} \cdot \operatorname{div} \cdot (\mathbf{v} - \mathbf{u}_0) \mathbf{d}\mathbf{x}. \end{aligned}$$

Then,

$$\bar{\mathbf{J}}(\mathbf{v}) = \sup_{\boldsymbol{\tau}_1 \in \boldsymbol{\Sigma}} \bar{\mathbf{L}}(\mathbf{v}; \boldsymbol{\tau}_1),$$

so that Problem $\bar{\mathcal{P}}$ is equivalent to the minimax problem

$\inf_{\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0} \sup_{\boldsymbol{\tau}_1 \in \boldsymbol{\Sigma}} \bar{\mathbf{L}}(\mathbf{v}; \boldsymbol{\tau}_1)$. The respective dual problem is

$$\sup_{\boldsymbol{\tau}_1 \in \boldsymbol{\Sigma}} \inf_{\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0} \bar{\mathbf{L}}(\mathbf{v}; \boldsymbol{\tau}_1).$$

Note that

$$\inf_{\mathbf{v} \in \mathbf{V}_0 + \mathbf{u}_0} \bar{\mathbf{L}}(\mathbf{v}; \boldsymbol{\tau}_1) = \begin{cases} \bar{\mathbf{I}}(\boldsymbol{\tau}_1) & \text{if } \boldsymbol{\tau}_1 \in \bar{\boldsymbol{\Sigma}}_{\mathbf{f}}(\Omega), \\ -\infty & \text{if } \boldsymbol{\tau}_1 \notin \bar{\boldsymbol{\Sigma}}_{\mathbf{f}}(\Omega), \end{cases}$$

where

$$\bar{\mathbf{I}}(\boldsymbol{\tau}_1) = \int_{\Omega} \left(\boldsymbol{\varepsilon}(\mathbf{u}_0) : (\boldsymbol{\tau}_1 + \boldsymbol{\tau}_2) - \frac{1}{2\nu} |\boldsymbol{\tau}_1|^2 - \psi^*(\boldsymbol{\tau}_2) - \mathbf{f} \cdot \mathbf{u}_0 \right) \mathbf{d}\mathbf{x},$$

and $\bar{\boldsymbol{\Sigma}}_{\mathbf{f}}$ is an affine subset in $\boldsymbol{\Sigma}$ that consists of the functions $\boldsymbol{\tau}$ satisfying (in a generalized sense) the condition $\mathbf{div}(\boldsymbol{\tau} + \boldsymbol{\tau}_2) = \nabla \mathbf{q} - \mathbf{f}$, i.e.,

$$\begin{aligned} \bar{\boldsymbol{\Sigma}}_{\mathbf{f}}(\Omega) &:= \left\{ \boldsymbol{\tau} \in \boldsymbol{\Sigma}(\Omega) \mid \int_{\Omega} \boldsymbol{\varepsilon}(\mathbf{w}) : (\boldsymbol{\tau} + \boldsymbol{\tau}_2) \mathbf{d}\mathbf{x} \right. \\ &= \left. \int_{\Omega} (\mathbf{f} \cdot \mathbf{w} + \mathbf{q} \operatorname{div} \mathbf{w}) \mathbf{d}\mathbf{x}, \mathbf{w} \in \mathbf{V}_0 \right\}. \end{aligned}$$

Thus, we arrive at the following formulation of the Dual Variational Problem.

Problem $\bar{\mathcal{P}}^*$ For given $\tau_2 \in \Sigma$ and $\mathbf{q} \in \mathring{\mathbf{L}}_2(\Omega)$ find a function $\bar{\sigma}_1 \in \bar{\Sigma}_f(\Omega)$ such that

$$\bar{\mathbf{I}}(\bar{\sigma}_1) = \sup \bar{\mathcal{P}}^* := \sup_{\tau_1 \in \bar{\Sigma}_f} \bar{\mathbf{I}}(\tau_1).$$

Theorem

Problem $\bar{\mathcal{P}}^$ has a unique solution $\bar{\sigma}_1$ satisfying the conditions*

$$\inf \bar{\mathcal{P}} = \bar{\mathbf{J}}(\bar{\mathbf{u}}) = \sup \bar{\mathcal{P}}^* = \bar{\mathbf{I}}(\bar{\sigma}_1), \quad (302)$$

$$\nu \varepsilon(\bar{\mathbf{u}}) = \bar{\sigma}_1. \quad (303)$$

Proof.

$-\bar{\mathbf{I}}$ is strictly convex and $\bar{\Sigma}_{\mathbf{f}}$ is a convex and closed subset of Σ . Therefore, Problem $\bar{\mathcal{P}}^*$ has a unique solution.

Note that $\sup \bar{\mathcal{P}}^* \leq \inf \bar{\mathcal{P}}$. This fact follows from the relation

$$\sup \inf \bar{\mathbf{L}} \leq \inf \sup \bar{\mathbf{L}}$$

$\bar{\mathbf{u}}$ satisfies the integral identity

$$\int_{\Omega} (\nu \varepsilon(\bar{\mathbf{u}}) : \varepsilon(\mathbf{w}) + \tau_2 : \varepsilon(\mathbf{w})) \, \mathbf{d}\mathbf{x} = \int_{\Omega} (\mathbf{f} \cdot \mathbf{w} + \mathbf{q} \operatorname{div} \mathbf{w}) \, \mathbf{d}\mathbf{x} \quad \mathbf{w} \in \mathbf{V}_0.$$

From here, it follows that

$$\begin{aligned} \int_{\Omega} \mathbf{f} \cdot \bar{\mathbf{u}} \, \mathbf{d}\mathbf{x} &= \int_{\Omega} (\nu |\varepsilon(\bar{\mathbf{u}})|^2 - \nu \varepsilon(\bar{\mathbf{u}}) : \varepsilon(\mathbf{u}_0) + \\ &+ \tau_2 : \varepsilon(\bar{\mathbf{u}} - \mathbf{u}_0) + \mathbf{f} \cdot \mathbf{u}_0 - \mathbf{q} \operatorname{div}(\bar{\mathbf{u}} - \mathbf{u}_0)) \, \mathbf{d}\mathbf{x}. \end{aligned}$$

Therefore,

$$\begin{aligned}
 \inf \bar{\mathcal{P}} &= \bar{\mathbf{J}}(\bar{\mathbf{u}}) \\
 &= \int_{\Omega} \left(\frac{\nu}{2} |\varepsilon(\bar{\mathbf{u}})|^2 + \boldsymbol{\tau}_2 : \varepsilon(\bar{\mathbf{u}}) - \psi^*(\boldsymbol{\tau}_2) \right) \mathbf{d}\mathbf{x} - \int_{\Omega} (\mathbf{q} \operatorname{div}(\bar{\mathbf{u}} - \mathbf{u}_0) + \mathbf{f} \cdot \bar{\mathbf{u}}) \mathbf{d}\mathbf{x} \\
 &= \int_{\Omega} \left(\boldsymbol{\tau}_2 : \varepsilon(\mathbf{u}_0) - \frac{\nu}{2} |\varepsilon(\bar{\mathbf{u}})|^2 - \psi^*(\boldsymbol{\tau}_2) + \nu \varepsilon(\mathbf{u}_0) : \varepsilon(\bar{\mathbf{u}}) \right) \mathbf{d}\mathbf{x} - \int_{\Omega} \mathbf{f} \cdot \mathbf{u}_0 \mathbf{d}\mathbf{x}.
 \end{aligned}$$

Since $\nu \varepsilon(\bar{\mathbf{u}}) \in \bar{\Sigma}_{\mathbf{f}}$, we know that

$$\bar{\mathbf{I}}(\nu \varepsilon(\bar{\mathbf{u}})) \leq \sup \bar{\mathcal{P}}^*$$

$$\begin{aligned}
\bar{\mathbf{I}}(\nu\varepsilon(\bar{\mathbf{u}})) &= \\
&= \int_{\Omega} \left(\varepsilon(\mathbf{u}_0) : (\nu\varepsilon(\bar{\mathbf{u}}) + \boldsymbol{\tau}_2) - \frac{\nu}{2} |\varepsilon(\bar{\mathbf{u}})|^2 - \psi^*(\boldsymbol{\tau}_2) \right) \mathbf{d}\mathbf{x} - \int_{\Omega} \mathbf{f} \cdot \mathbf{u}_0 \mathbf{d}\mathbf{x} \leq \\
&\leq \sup \bar{\mathcal{P}}^* \leq \inf \bar{\mathcal{P}} = \\
&= \int_{\Omega} \left(\boldsymbol{\tau}_2 : \varepsilon(\mathbf{u}_0) - \frac{\nu}{2} |\varepsilon(\bar{\mathbf{u}})|^2 - \psi^*(\boldsymbol{\tau}_2) + \nu\varepsilon(\mathbf{u}_0) : \varepsilon(\bar{\mathbf{u}}) \right) \mathbf{d}\mathbf{x} - \int_{\Omega} \mathbf{f} \cdot \mathbf{u}_0 \mathbf{d}\mathbf{x}.
\end{aligned}$$

Consequently, $\nu\varepsilon(\bar{\mathbf{u}}) = \bar{\boldsymbol{\sigma}}_1$ and

$$\inf \bar{\mathcal{P}} = \sup \bar{\mathcal{P}}^*.$$

Estimates of deviations from exact solutions for solenoidal fields

Let $\mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0$. For it and $\boldsymbol{\tau}_{\mathbf{f}} \in \bar{\boldsymbol{\Sigma}}_{\mathbf{f}}(\Omega)$ the inequality

$$\begin{aligned} \int_{\Omega} \frac{\nu}{2} |\boldsymbol{\varepsilon}(\mathbf{v} - \mathbf{u})|^2 \, \mathbf{d}\mathbf{x} &\leq \\ &\leq \mathbf{J}(\mathbf{v}) - \inf \bar{\mathcal{P}} = \mathbf{J}(\mathbf{v}) - \sup \bar{\mathcal{P}}^* \leq \mathbf{J}(\mathbf{v}) - \bar{\mathbf{I}}(\boldsymbol{\tau}_{\mathbf{f}}) \end{aligned} \quad (304)$$

holds. Estimate the right-hand side of (304) as follows:

$$\begin{aligned} \mathbf{J}(\bar{\mathbf{v}}) - \bar{\mathbf{I}}(\boldsymbol{\tau}_{\mathbf{f}}) &\leq \int_{\Omega} \left(\frac{\nu}{2} |\boldsymbol{\varepsilon}(\mathbf{v})|^2 + \frac{1}{2\nu} |\boldsymbol{\tau}_{\mathbf{f}}|^2 - \boldsymbol{\varepsilon}(\mathbf{u}_0) : \boldsymbol{\tau}_{\mathbf{f}} \right) \, \mathbf{d}\mathbf{x} \\ &+ \int_{\Omega} (\psi(\boldsymbol{\varepsilon}(\mathbf{v})) + \psi^*(\boldsymbol{\tau}_2) - \boldsymbol{\varepsilon}(\mathbf{u}_0) : \boldsymbol{\tau}_2) \, \mathbf{d}\mathbf{x} + \int_{\Omega} \mathbf{f} \cdot (\mathbf{u}_0 - \mathbf{v}) \, \mathbf{d}\mathbf{x}. \end{aligned}$$

Let $\mathbf{q} \in \mathring{\mathbf{L}}_2(\Omega)$. Since $\boldsymbol{\tau}_{\mathbf{f}} \in \bar{\boldsymbol{\Sigma}}_{\mathbf{f}}(\Omega)$ and $\mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0$, we have

$$\begin{aligned} \int_{\Omega} \mathbf{f} \cdot (\mathbf{u}_0 - \mathbf{v}) \, d\mathbf{x} &= \int_{\Omega} (\mathbf{f} \cdot (\mathbf{u}_0 - \mathbf{v}) + \mathbf{q} \operatorname{div}(\mathbf{u}_0 - \mathbf{v})) \, d\mathbf{x} \\ &= \int_{\Omega} \boldsymbol{\varepsilon}(\mathbf{u}_0 - \mathbf{v}) : (\boldsymbol{\tau}_{\mathbf{f}} + \boldsymbol{\tau}_2) \, d\mathbf{x}. \end{aligned}$$

As a result, we obtain the following estimate:

$$\frac{\nu}{2} \|\boldsymbol{\varepsilon}(\mathbf{v} - \mathbf{u})\|^2 \leq \mathcal{M}_1(\mathbf{v}, \boldsymbol{\tau}_{\mathbf{f}}, \boldsymbol{\tau}_2) := \mathbf{D}_1(\boldsymbol{\varepsilon}(\mathbf{v}), \boldsymbol{\tau}_{\mathbf{f}}) + \mathbf{D}_2(\boldsymbol{\varepsilon}(\mathbf{v}), \boldsymbol{\tau}_2), \quad (305)$$

Here we have introduced *compound functionals*

$$\begin{aligned} \mathbf{D}_1(\boldsymbol{\varepsilon}(\mathbf{v}), \boldsymbol{\tau}_{\mathbf{f}}) &:= \int_{\Omega} \left(\frac{\nu}{2} |\boldsymbol{\varepsilon}(\mathbf{v})|^2 + \frac{1}{2\nu} |\boldsymbol{\tau}_{\mathbf{f}}|^2 - \boldsymbol{\varepsilon}(\mathbf{v}) : \boldsymbol{\tau}_{\mathbf{f}} \right) \mathbf{d}\mathbf{x} = \\ &= \frac{1}{2\nu} \|\nu\boldsymbol{\varepsilon}(\mathbf{v}) - \boldsymbol{\tau}_{\mathbf{f}}\|^2; \\ \mathbf{D}_2(\boldsymbol{\varepsilon}(\mathbf{v}), \boldsymbol{\tau}_2) &:= \int_{\Omega} (\psi(\boldsymbol{\varepsilon}(\mathbf{v})) + \psi^*(\boldsymbol{\tau}_2) - \boldsymbol{\varepsilon}(\mathbf{v}) : \boldsymbol{\tau}_2) \mathbf{d}\mathbf{x}. \end{aligned}$$

It is clear that both functionals \mathbf{D}_1 and \mathbf{D}_2 are nonnegative. Moreover,

$$\mathbf{D}_1(\boldsymbol{\varepsilon}(\mathbf{v}), \boldsymbol{\tau}_{\mathbf{f}}) = 0$$

if and only if

$$\boldsymbol{\tau}_{\mathbf{f}} = \nu\boldsymbol{\varepsilon}(\mathbf{v}).$$

By the properties of the conjugate functionals

$$\mathbf{D}_2(\boldsymbol{\varepsilon}(\mathbf{v}), \boldsymbol{\tau}_2) = \mathbf{0}$$

if and only if

$$\boldsymbol{\tau}_2 \in \partial\psi(\boldsymbol{\varepsilon}(\mathbf{v})).$$

Now, it is easy to understand the meaning of the estimates (305). Present the main system in the form

$$-\mathbf{div}(\boldsymbol{\sigma}_1 + \boldsymbol{\sigma}_2) = \mathbf{f} - \nabla \mathbf{p} \quad \text{in } \Omega, \quad (306)$$

$$\mathbf{div} \mathbf{u} = \mathbf{0} \quad \text{in } \Omega, \quad (307)$$

$$\boldsymbol{\sigma}_1 = \nu \boldsymbol{\varepsilon}(\mathbf{u}), \boldsymbol{\sigma}_2 \in \partial\psi(\boldsymbol{\varepsilon}(\mathbf{u})) \quad \text{in } \Omega, \quad (308)$$

$$\mathbf{u} = \mathbf{u}_0 \quad \text{on } \partial\Omega. \quad (309)$$

If $\mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0$, and $\boldsymbol{\tau}_{\mathbf{f}} \in \bar{\boldsymbol{\Sigma}}_{\mathbf{f}}(\Omega)$, then

$$-\mathbf{div}(\boldsymbol{\tau}_{\mathbf{f}} + \boldsymbol{\tau}_2) = \mathbf{f} - \nabla \mathbf{q} \text{ and } \mathbf{div} \mathbf{v} = \mathbf{0},$$

so that for \mathbf{v} , $\boldsymbol{\tau}_{\mathbf{f}}$, $\boldsymbol{\tau}_2$ and q the relations (306), (307) and (309) are satisfied.

Our estimate shows that in this case, the energy norm of the deviation from the exact solution is controlled by the quantities

$$\mathbf{D}_1(\boldsymbol{\varepsilon}(\mathbf{v}), \boldsymbol{\tau}_{\mathbf{f}}) \text{ and } \mathbf{D}_2(\boldsymbol{\varepsilon}(\mathbf{v}), \boldsymbol{\tau}_2),$$

Certainly, the condition $\boldsymbol{\tau}_f \in \bar{\boldsymbol{\Sigma}}_f(\Omega)$ is rather obligatory and it would be useful to somehow eliminate it. This can be done by the same method as we have discussed for linear problems.

As a result, we obtain the estimate

$$\begin{aligned} \frac{\nu}{2} \|\boldsymbol{\varepsilon}(\mathbf{v} - \mathbf{u})\|^2 &\leq (1 + \beta) \mathbf{D}_1(\boldsymbol{\varepsilon}(\mathbf{v}), \boldsymbol{\tau}_1) + \mathbf{D}_2(\boldsymbol{\varepsilon}(\mathbf{v}), \boldsymbol{\tau}_2) \\ &\quad + \frac{1 + \beta}{2\nu\beta} \|\mathbf{div}(\boldsymbol{\tau}_1 + \boldsymbol{\tau}_2) + \mathbf{f} - \nabla \mathbf{q}\|^2, \end{aligned} \quad (310)$$

that holds for any function $\mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0$, any pair of functions $(\boldsymbol{\tau}_1, \boldsymbol{\tau}_2) \in \boldsymbol{\Sigma} \times \boldsymbol{\Sigma}$, and any $\beta > 0$. The right-hand side of (310) is the majorant of the norm of the deviation from the exact solution that we denote by $\mathcal{M}_2(\beta, \mathbf{v}, \boldsymbol{\tau}_1, \boldsymbol{\tau}_2, \mathbf{q})$.

If the sum $\boldsymbol{\tau}_1 + \boldsymbol{\tau}_2$ has a higher regularity, so that

$$(\boldsymbol{\tau}_1 + \boldsymbol{\tau}_2) \in \boldsymbol{\Sigma}_{\text{div}}(\Omega)$$

and, in addition, $\mathbf{q} \in \mathbf{H}^1$, then the last term of the Majorant is estimated by an explicitly computable integral:

$$\begin{aligned} \frac{\nu}{2} \|\boldsymbol{\varepsilon}(\mathbf{v} - \mathbf{u})\|^2 \leq & (\mathbf{1} + \beta) \mathbf{D}_1(\boldsymbol{\varepsilon}(\mathbf{v}), \boldsymbol{\tau}_1) + \mathbf{D}_2(\boldsymbol{\varepsilon}(\mathbf{v}), \boldsymbol{\tau}_2) + \\ & + \frac{\mathbf{1} + \beta}{2\nu\beta} \mathbf{C}_\Omega^2 \|\text{div}(\boldsymbol{\tau}_1 + \boldsymbol{\tau}_2) + \mathbf{f} - \nabla \mathbf{q}\|^2. \quad (311) \end{aligned}$$

Theorem

For any $\beta \in \mathbb{R}_+$, $\mathbf{v} \in \mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0$, $\boldsymbol{\tau}_1 \in \boldsymbol{\Sigma}$, $\boldsymbol{\tau}_2 \in \boldsymbol{\Sigma}$ and $\mathbf{q} \in \mathring{\mathbf{L}}_2(\Omega)$ the functional $\mathcal{M}_2(\beta, \mathbf{v}, \boldsymbol{\tau}_1, \boldsymbol{\tau}_2, \mathbf{q})$ majorizes the quantity $\|\varepsilon(\mathbf{v} - \mathbf{u})\|^2$.

For any $\beta \in \mathbb{R}_+$, infimum of this functional on the set

$$(\mathring{\mathbf{J}}_2^1(\Omega) + \mathbf{u}_0) \times \boldsymbol{\Sigma} \times \boldsymbol{\Sigma} \times \mathring{\mathbf{L}}_2(\Omega)$$

is equal to zero and it is attained if and only if $\mathbf{v} = \mathbf{u}$, $\boldsymbol{\tau}_1 = \boldsymbol{\sigma}_1$, $\boldsymbol{\tau}_2 = \boldsymbol{\sigma}_2$ and $\mathbf{q} = \mathbf{p}$.

Example 1.

For the Stokes problem $\mathbf{K}_* = \mathbf{0}$, $\psi(\boldsymbol{\varepsilon}) \equiv \mathbf{0}$. Set $\boldsymbol{\tau}_2 = \mathbf{0}$. Then, $\mathbf{D}_2(\boldsymbol{\varepsilon}(\boldsymbol{\nu}), \boldsymbol{\tau}_2) \equiv 0$ and (310) comes in the form we have already obtained.

Example 2.

For the Bingham model $\psi(\boldsymbol{\varepsilon}) = \mathbf{K}_* |\boldsymbol{\varepsilon}|$, and

$$\psi^*(\boldsymbol{\tau}(\mathbf{x})) = \begin{cases} \mathbf{0}, & \text{if } |\boldsymbol{\tau}(\mathbf{x})| \leq \mathbf{K}_*, \\ +\infty, & \text{if } |\boldsymbol{\tau}(\mathbf{x})| > \mathbf{K}_*. \end{cases}$$

Therefore,

$$\mathbf{D}_2(\boldsymbol{\varepsilon}(\mathbf{v}), \boldsymbol{\tau}_2) = \int_{\Omega} (\mathbf{k}_* |\boldsymbol{\varepsilon}(\mathbf{v})| - \boldsymbol{\varepsilon}(\mathbf{v}) \cdot \boldsymbol{\tau}_2) \mathbf{d}\mathbf{x},$$

if almost everywhere the function $\boldsymbol{\tau}_2$ satisfies the condition $|\boldsymbol{\tau}_2(\mathbf{x})| \leq \mathbf{K}_*$. In the opposite case, $\mathbf{D}_2(\boldsymbol{\varepsilon}(\mathbf{v}), \boldsymbol{\tau}_2) = +\infty$. Then, (310) comes in the form:

$$\begin{aligned} \frac{\nu}{2} \|\boldsymbol{\varepsilon}(\mathbf{v} - \mathbf{u})\|^2 &\leq \int_{\Omega} \left(\frac{(1 + \beta)}{2\nu} |\nu\boldsymbol{\varepsilon}(\mathbf{v}) - \boldsymbol{\tau}_1|^2 + \mathbf{K}_* |\boldsymbol{\varepsilon}(\mathbf{v})| - \boldsymbol{\varepsilon}(\mathbf{v}) : \boldsymbol{\tau}_2 \right) \mathbf{d}\mathbf{x} \\ &+ \frac{1 + \beta}{2\beta\nu} \mathbf{I} \mathbf{div}(\boldsymbol{\tau}_1 + \boldsymbol{\tau}_2) + \mathbf{f} - \nabla \mathbf{q} \mathbf{I}^2. \end{aligned} \quad (312)$$

If the functions $\boldsymbol{\tau}_1$, $\boldsymbol{\tau}_2$ and \mathbf{q} are taken such that $\mathbf{q} \in \mathbf{H}^1$ $\mathbf{div}(\boldsymbol{\tau}_1 + \boldsymbol{\tau}_2) \in \mathbf{L}_2(\Omega)$, then the last term in (312) is estimated by the integral.

Bingham fluid may have two zones: the *congestion zone* Ω_0 (where $\boldsymbol{\varepsilon}(\mathbf{u}) \equiv \mathbf{0}$) and the *flow zone* Ω_1 (where $|\boldsymbol{\varepsilon}(\mathbf{u})| > 0$).

Assume that the right-hand side of (312) vanishes for certain functions \mathbf{v} , $\boldsymbol{\tau}_1$, $\boldsymbol{\tau}_2$ and \mathbf{q} . Then, in Ω_0 we have $\boldsymbol{\varepsilon}(\mathbf{v}) = \mathbf{0}$, and, consequently, $\boldsymbol{\tau}_1 = \mathbf{0}$ and $\mathbf{div}\boldsymbol{\tau}_2 + \mathbf{f} - \nabla\mathbf{q} = \mathbf{0}$ for some $\boldsymbol{\tau}_2$, satisfying the condition $|\boldsymbol{\tau}_2(\mathbf{x})| \leq \mathbf{1}$.

At the same time, in the flow zone Ω_1 the relations

$$\boldsymbol{\tau}_2 = \mathbf{k}_* \frac{\boldsymbol{\varepsilon}(\mathbf{v})}{|\boldsymbol{\varepsilon}(\mathbf{v})|}, \quad \boldsymbol{\tau}_1 = \nu\boldsymbol{\varepsilon}(\mathbf{v}), \quad \mathbf{div}(\boldsymbol{\tau}_1 + \boldsymbol{\tau}_2) + \mathbf{f} - \nabla\mathbf{q} = \mathbf{0}$$

hold.

Lecture 9

GENERAL APPROACH TO A POSTERIORI ERROR CONTROL FOR NONLINEAR VARIATIONAL PROBLEMS

Lecture goal

In subsequent lectures we will present the general theory of a posteriori error control for convex variational problems. In the framework of this theory we are able to derive computable upper bounds for the errors for problems of the type

$$\inf_{\mathbf{v} \in \mathbf{V}} \mathbf{J}(\mathbf{v}, \Lambda \mathbf{v}), \quad \mathbf{J}(\mathbf{v}, \Lambda \mathbf{v}) := \mathbf{G}(\Lambda \mathbf{v}) + \mathbf{F}(\mathbf{v}),$$

where $\Lambda : \mathbf{V} \rightarrow \mathbf{Y}$ is a linear continuous operator from a Banach space \mathbf{V} to another Banach space \mathbf{Y} and $\mathbf{J} : \mathbf{Y} \rightarrow \mathbb{R}$ and $\mathbf{F} : \mathbf{V} \rightarrow \mathbb{R}$ are convex l.s.c. functionals.

In particular, if

$$\Lambda \mathbf{v} = \nabla \mathbf{v}, \quad \mathbf{G}(\mathbf{y}) = (\mathbf{A}\mathbf{y}, \mathbf{y}), \quad \mathbf{F}(\mathbf{v}) = (\mathbf{f}, \mathbf{v}),$$

then we arrive to the variational formulation of the problem

$$\mathbf{div} \mathbf{A} \nabla \mathbf{u} + \mathbf{f} = \mathbf{0}.$$

Many other problems have the above form, were

G is the **energy functional** whose form is dictated by the **dissipative properties of a media**.

F is the functional associated with **external forces** and (or) **boundary conditions**.

Diffusion type problems,
Linear elasticity,
Biharmonic problems,
Kirghoff and Mindlin plates,
Problems in deformation theory of elastoplasticity,
p-Laplace equation,
Stokes problem,
Nonlinear problems in the theory of viscous fluids and many other problems can be presented in the above general form.

In such models, the structure of the "energy functional" G plays crucial role in all the parts of the mathematical analysis: existence and differentiability properties of minimizers and estimates of deviations from the minimizers.

To understand the basic principles of the functional approach to the derivation of a posteriori bounds of the approximation errors we need to make a **concise overview of some parts of the duality theory in the calculus of variations.**

A consequent exposition functional type a posteriori error estimates for nonlinear variational problems can be found in the papers

S. Repin. A posteriori error estimation for variational problems with uniformly convex functionals, *Math. Comput.*, 69(230), 2000, 481-500.

S. Repin. Two-sided estimates for deviation from an exact solution to uniformly elliptic equation. *Trudi St.-Petersburg Math. Society*, 9(2001), 148-179 (in Russian, translated in *American Mathematical Translations Series 2*, 9(2003)

and in the book

P. Neittaanmaki and S. Repin. Reliable methods for computer simulation. Error control and a posteriori estimates. Elsevier, NY, 2004.

Selected topics of the duality theory in the calculus of variations

To understand the structure of functional a posteriori estimates for the considered class of problems we need first discuss three additional topics:

- **Dual and bidual functionals ;**
- **Compound functionals ;**
- **Uniformly convex functionals.**

Dual (polar) functionals

Hereafter \mathbf{V}^* contains all linear continuous functionals defined on \mathbf{V} . The elements of \mathbf{V}^* are marked by stars,

$\langle \mathbf{v}^*, \mathbf{v} \rangle$ is called the **duality pairing** of the spaces \mathbf{V} and \mathbf{V}^* .

Let $\mathbf{J} : \mathbf{V} \rightarrow \mathbb{R}$, then \mathbf{J}^* defined by the relation

$$\mathbf{J}^*(\mathbf{v}^*) = \sup_{\mathbf{v} \in \mathbf{V}} \{ \langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{J}(\mathbf{v}) \}$$

is called **dual** to \mathbf{J} .

If \mathbf{J} is a smooth function that increases at infinity faster than any linear function, then \mathbf{J}^* is the Legendre transform of \mathbf{J} . The above general definition comes from Young and Fenchel. The functional \mathbf{J}^* is also called **polar** to \mathbf{J} .

Bipolar functionals

The functional

$$\mathbf{J}^{**}(\mathbf{v}) = \sup_{\mathbf{v}^* \in \mathbf{V}^*} \{\langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{J}^*(\mathbf{v}^*)\}$$

is called the **bidual** to \mathbf{J} (or **bipolar**).

Straightforwardly from the definition, it follows that \mathbf{J}^* and \mathbf{J}^{**} are convex functionals (they are defined as upper bounds of affine functionals).

Formally, one can also define

$$\mathbf{J}^{***}(\mathbf{v}^*) := \sup_{\mathbf{v} \in \mathbf{V}} \{\langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{J}^{**}(\mathbf{v})\}.$$

However, this definition brings nothing new. It is proved that

$$\mathbf{J}^{***}(\mathbf{v}^*) = \mathbf{J}^*(\mathbf{v}^*), \quad \forall \mathbf{v}^* \in \mathbf{V}^*.$$

Mutually dual functionals

Let $\mathbf{J} : \mathbf{V} \rightarrow \overline{\mathbb{R}} := \{\mathbb{R}, -\infty, +\infty\}$ and $\mathbf{G}^* : \mathbf{V}^* \rightarrow \overline{\mathbb{R}}$ be two functionals defined on a Banach space \mathbf{V} and its dual space \mathbf{V}^* , respectively. These two functionals are called **mutually dual** if

$$(\mathbf{G}^*)^* = \mathbf{J} \quad \text{and} \quad \mathbf{J}^* = \mathbf{G}^*.$$

Examples

To illustrate the definitions of conjugate functionals, we present below several examples for functionals defined on the Euclidean space \mathbf{E}^d . In this case, \mathbf{V} and \mathbf{V}^* are isometrically isomorphic. Their elements are d -dimensional vectors denoted by $\boldsymbol{\xi}$ and $\boldsymbol{\xi}^*$, respectively, so that

$$\langle \boldsymbol{\xi}^*, \boldsymbol{\xi} \rangle = \boldsymbol{\xi}^* \cdot \boldsymbol{\xi} = \xi_i^* \xi_i.$$

These examples have a practical meaning because for a wide class of integral type functionals (in the mechanics they are the **energy functionals**) finding the dual energy functional is reduced to **finding dual to its integrand !**

In other words, if the "primal energy functional" has the form

$$\mathbf{G}(\mathbf{v}) := \int_{\Omega} \mathbf{g}(\Lambda \mathbf{v}) \mathbf{d}\mathbf{x}$$

where \mathbf{g} is the "internal energy" or "dissipative potential", then the so-called "complementary energy" is given by the integral functional

$$\mathbf{G}^*(\mathbf{y}^*) := \int_{\Omega} \mathbf{g}^*(\mathbf{y}^*) \mathbf{d}\mathbf{x},$$

where \mathbf{g}^* is conjugate to \mathbf{g} in the algebraic sense.

Example 1 (Diffusion problems)

Let $\mathbf{A} = \{\mathbf{a}_{ij}\}$ be a real, positive definite matrix and

$$\mathbf{g}(\boldsymbol{\xi}) = \frac{1}{2}\mathbf{A}\boldsymbol{\xi} \cdot \boldsymbol{\xi} = \frac{1}{2}\mathbf{a}_{ij}\xi_i\xi_j.$$

Then

$$\mathbf{g}^*(\boldsymbol{\xi}^*) = \sup_{\boldsymbol{\xi} \in \mathbb{E}^d} \left\{ \boldsymbol{\xi}^* \cdot \boldsymbol{\xi} - \frac{1}{2}\mathbf{A}\boldsymbol{\xi} \cdot \boldsymbol{\xi} \right\}.$$

This supremum is attained on an element $\boldsymbol{\xi}_0$ such that

$$\boldsymbol{\xi}^* = \mathbf{A}\boldsymbol{\xi}_0 \implies \boldsymbol{\xi}_0 = \mathbf{A}^{-1}\boldsymbol{\xi}^*.$$

Therefore, we have a pair of mutually conjugate functionals

$$\mathbf{g}(\boldsymbol{\xi}) = \frac{1}{2}\mathbf{A}\boldsymbol{\xi} \cdot \boldsymbol{\xi} \quad \text{and} \quad \mathbf{g}^*(\boldsymbol{\xi}^*) = \frac{1}{2}\mathbf{A}^{-1}\boldsymbol{\xi}^* \cdot \boldsymbol{\xi}^*.$$

In diffusion type boundary-value problems we arrive at the functional (with $\mathbf{y} = \nabla \mathbf{v}$)

$$\frac{1}{2} \int_{\Omega} \mathbf{A} \mathbf{y} \cdot \mathbf{y} \, dx \quad \mathbf{y} \in \mathbf{L}^2(\Omega, \mathbb{R}^n),$$

which is mutually dual to

$$\frac{1}{2} \int_{\Omega} \mathbf{A}^{-1} \mathbf{y}^* \cdot \mathbf{y}^* \, dx \quad \mathbf{y}^* \in \mathbf{L}^2(\Omega, \mathbb{R}^n)$$

Example 2 (Linear elasticity)

Let $\mathbf{L} = \{\mathbf{L}_{ijkl}\}$ be a real, positive definite tensor of the 4-th order and $\boldsymbol{\tau}$ be a tensor of the second order ($\mathbf{d} \times \mathbf{d}$ -matrix). Then,

$$\mathbf{g}(\boldsymbol{\varepsilon}) = \frac{1}{2} \mathbf{L} \boldsymbol{\varepsilon} : \boldsymbol{\varepsilon} = \frac{1}{2} \mathbf{L}_{ijkl} \varepsilon_{ij} \varepsilon_{km}.$$

Then

$$\mathbf{g}^*(\boldsymbol{\varepsilon}^*) = \sup_{\boldsymbol{\varepsilon} \in \mathbf{M}^{\mathbf{d} \times \mathbf{d}}} \left\{ \boldsymbol{\varepsilon}^* : \boldsymbol{\varepsilon} - \frac{1}{2} \mathbf{A} \boldsymbol{\varepsilon} : \boldsymbol{\varepsilon} \right\}.$$

This supremum is attained on an element $\boldsymbol{\varepsilon}_0$ such that

$$\boldsymbol{\tau}^* = \mathbf{L} \boldsymbol{\varepsilon}_0 \implies \boldsymbol{\varepsilon}_0 = \mathbf{L}^{-1} \boldsymbol{\varepsilon}^*.$$

Therefore, we have a pair of mutually dual functionals

$$\mathbf{g}(\boldsymbol{\varepsilon}) = \frac{1}{2} \mathbf{L} \boldsymbol{\varepsilon} : \boldsymbol{\varepsilon} \quad \text{and} \quad \mathbf{g}^*(\boldsymbol{\varepsilon}^*) = \frac{1}{2} \mathbf{L}^{-1} \boldsymbol{\varepsilon}^* : \boldsymbol{\varepsilon}^*.$$

In linear elasticity problems we arrive at the energy functional in terms of strains $\varepsilon(\mathbf{v}) = \frac{1}{2}(\nabla\mathbf{v} + (\nabla\mathbf{v})^T)$

$$\frac{1}{2} \int_{\Omega} \mathbb{L}\varepsilon : \varepsilon \, dx \quad \varepsilon \in \mathbf{L}^2(\Omega, \mathbb{M}^{n \times n}),$$

which is mutually dual to the "complementary energy" functional written in terms of stresses $\varepsilon^*(x) \Rightarrow \tau(x)$

$$\frac{1}{2} \int_{\Omega} \mathbb{L}^{-1}\tau : \tau \, dx \quad \tau \in \mathbf{L}^2(\Omega, \mathbb{M}^{n \times n})$$

Example 3 (Nonlinear elasticity, p-Laplacian)

Consider the functional

$$\mathbf{g}(\boldsymbol{\xi}) = \frac{1}{\mathbf{p}} |\boldsymbol{\xi}|^{\mathbf{p}},$$

where $\mathbf{p} > 1$ and $|\boldsymbol{\xi}| = (\boldsymbol{\xi} \cdot \boldsymbol{\xi})^{1/2}$. It is easy to verify that the quantity $\boldsymbol{\xi}^* \cdot \boldsymbol{\xi} - \frac{1}{\mathbf{p}} |\boldsymbol{\xi}|^{\mathbf{p}}$ attains a supremum if $\boldsymbol{\xi} = \boldsymbol{\xi}_0$, where $\boldsymbol{\xi}_0$ satisfies the relation

$$\boldsymbol{\xi}^* - |\boldsymbol{\xi}_0|^{\mathbf{p}-2} \boldsymbol{\xi}_0 = \mathbf{0},$$

which yields $|\boldsymbol{\xi}^*| = |\boldsymbol{\xi}_0|^{\mathbf{p}-1}$ and $\boldsymbol{\xi}^* \cdot \boldsymbol{\xi}_0 = |\boldsymbol{\xi}_0|^{\mathbf{p}}$. Therefore,

$$\mathbf{g}^*(\boldsymbol{\xi}^*) = \boldsymbol{\xi}^* \cdot \boldsymbol{\xi}_0 - \frac{1}{\mathbf{p}} |\boldsymbol{\xi}_0|^{\mathbf{p}} = \left(1 - \frac{1}{\mathbf{p}}\right) |\boldsymbol{\xi}_0|^{\mathbf{p}} = \frac{1}{\mathbf{p}^*} |\boldsymbol{\xi}^*|^{\mathbf{p}^*},$$

where $\mathbf{p}^* = \frac{\mathbf{p}}{\mathbf{p}-1}$.

Thus, we obtain another pair of mutually conjugate functionals

$$\mathbf{g}(\xi) = \frac{1}{\mathbf{p}} |\xi|^{\mathbf{p}} \quad \text{and} \quad \mathbf{g}^*(\xi^*) = \frac{1}{\mathbf{p}^*} |\xi^*|^{\mathbf{p}^*},$$

$$\text{where } \frac{1}{\mathbf{p}} + \frac{1}{\mathbf{p}^*} = 1.$$

Remark

This relation admits generalizations. Namely, let $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ be a proper convex function that is, in addition, odd and let $\varphi^* : \mathbb{R} \rightarrow \mathbb{R}$ be its conjugate. Then

$$(\varphi(\|\mathbf{u}\|\mathbf{v}))^* = \varphi^*(\|\mathbf{u}^*\|\mathbf{v}^*).$$

In certain nonlinear boundary-value problems we arrive at the functional
(with $\mathbf{y} = \nabla \mathbf{v}$ or $\mathbf{y} = \boldsymbol{\varepsilon}(\mathbf{v})$)

$$\frac{1}{p} \int_{\Omega} |\mathbf{y}|^p \, dx \quad \mathbf{y} \in L^p(\Omega, \mathbb{R}^n[M^{n \times n}]),$$

which is mutually dual to

$$\frac{1}{p^*} \int_{\Omega} |\mathbf{y}^*|^{p^*} \, dx \quad \mathbf{y}^* \in L^{p^*}(\Omega, \mathbb{R}^n[M^{n \times n}]).$$

Example 4 (Action of external forces)

Let $\mathbf{g}(\boldsymbol{\xi})$ be a linear functional, i.e.,

$$\mathbf{g}(\boldsymbol{\xi}) = \ell \cdot \boldsymbol{\xi}, \quad \ell \in \mathbf{E}^d.$$

It is easy to see that

$$\mathbf{g}^*(\boldsymbol{\xi}^*) = \sup_{\boldsymbol{\xi} \in \mathbf{E}^d} \{\boldsymbol{\xi}^* \cdot \boldsymbol{\xi} - \ell \cdot \boldsymbol{\xi}\} = \begin{cases} 0 & \boldsymbol{\xi}^* = \ell, \\ +\infty & \boldsymbol{\xi}^* \neq \ell. \end{cases}$$

Denote by $\mathfrak{X}_{\{\ell\}}$ the characteristic functional of the set $\{\ell\} \subset \mathbf{E}^d$. Then, another pair of mutually conjugate functionals is as follows:

$$\mathbf{g}(\boldsymbol{\xi}) = \ell \cdot \boldsymbol{\xi} \quad \text{and} \quad \mathbf{g}^*(\boldsymbol{\xi}^*) = \mathfrak{X}_{\{\ell\}}(\boldsymbol{\xi}^*).$$

Thus, for the functional $\mathbf{G} : \mathbf{L}^2 \rightarrow \mathbb{R}$

$$\mathbf{G}(\mathbf{v}) := \int_{\Omega} \mathbf{f} \mathbf{v} \, dx, \quad \mathbf{f} \in \mathbf{L}^2(\Omega)$$

the respective dual functional is $\mathbf{G}^* : \mathbf{L}^2 \rightarrow \mathbb{R}$

$$\mathbf{G}^*(\mathbf{v}^*) = \mathbf{0} \text{ if } \mathbf{v}^* = \mathbf{f}, \quad \mathbf{G}^*(\mathbf{v}^*) = +\infty \text{ in other cases.}$$

Example 5 (Variational inequalities)

Let $\mathbf{g}(\boldsymbol{\xi}) = |\boldsymbol{\xi}|$. Then

$$\sup_{\boldsymbol{\xi}} \{\boldsymbol{\xi}^* \cdot \boldsymbol{\xi} - |\boldsymbol{\xi}|\}$$

may be finite or infinite depending on the value of $|\boldsymbol{\xi}^*|$. If $|\boldsymbol{\xi}^*| > 1$, then, obviously, it is infinite. If $|\boldsymbol{\xi}^*| \leq 1$, then, on the one hand,

$$\sup_{\boldsymbol{\xi}} \{\boldsymbol{\xi}^* \cdot \boldsymbol{\xi} - |\boldsymbol{\xi}|\} \leq \sup_{\boldsymbol{\xi}} \{1|\boldsymbol{\xi}| - |\boldsymbol{\xi}|\} = 0.$$

On the other hand, $\sup_{\boldsymbol{\xi}} \{\boldsymbol{\xi}^* \cdot \boldsymbol{\xi} - |\boldsymbol{\xi}|\} \geq \boldsymbol{\xi}^* \cdot \mathbf{0} - 0 = 0$. This means that $\mathbf{g}^*(\boldsymbol{\xi}^*) = \mathbf{0}$ if $|\boldsymbol{\xi}^*| \leq 1$ and, thus,

$$\mathbf{g}(\boldsymbol{\xi}) = |\boldsymbol{\xi}|, \quad \mathbf{g}^*(\boldsymbol{\xi}^*) = \mathfrak{X}_{\mathcal{B}^*(\mathbf{0},1)}(\boldsymbol{\xi}^*), \quad \text{where } \mathcal{B}^*(\mathbf{0},1) = \{\boldsymbol{\xi}^* \in \mathbf{E}^d \mid |\boldsymbol{\xi}^*| \leq 1\}.$$

Thus, for the functional $\mathbf{G} : \mathbf{L}^1 \rightarrow \mathbb{R}$

$$\mathbf{G}(\mathbf{v}) := \int_{\Omega} |\mathbf{v}| \, d\mathbf{x},$$

the respective dual functional is $\mathbf{G}^* : \mathbf{L}^\infty \rightarrow \mathbb{R}$

$$\mathbf{G}^*(\mathbf{v}^*) = 0 \text{ if } |\mathbf{v}^*(\mathbf{x})| \leq 1 \text{ a.e. in } \Omega, \quad \mathbf{G}^*(\mathbf{v}^*) = +\infty \text{ in other cases.}$$

Properties of dual functionals

Property 1

If $\mathbf{J} : \mathbf{V} \rightarrow \overline{\mathbb{R}}$ and $\mathbf{G} : \mathbf{V} \rightarrow \overline{\mathbb{R}}$ are such that

$$\mathbf{J}(\mathbf{v}) \geq \mathbf{G}(\mathbf{v}), \quad \forall \mathbf{v} \in \mathbf{V},$$

then

$$\mathbf{J}^*(\mathbf{v}^*) \leq \mathbf{G}^*(\mathbf{v}^*), \quad \forall \mathbf{v}^* \in \mathbf{V}^*.$$

Proof. We have

$$\mathbf{J}^*(\mathbf{v}^*) = \sup_{\mathbf{v} \in \mathbf{V}} \{\langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{J}(\mathbf{v})\} \leq \sup_{\mathbf{v} \in \mathbf{V}} \{\langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{G}(\mathbf{v})\} = \mathbf{G}^*(\mathbf{v}^*).$$

Property 2

For any $\lambda > 0$,

$$(\lambda \mathbf{J})^*(\mathbf{v}^*) = \lambda \mathbf{J}^*\left(\frac{\mathbf{v}^*}{\lambda}\right).$$

Proof. This property is justified by direct calculations:

$$\begin{aligned} (\lambda \mathbf{J})^*(\mathbf{v}^*) &= \sup_{\mathbf{v} \in \mathbf{V}} \{ \langle \mathbf{v}^*, \mathbf{v} \rangle - \lambda \mathbf{J}(\mathbf{v}) \} = \\ &= \lambda \sup_{\mathbf{v} \in \mathbf{V}} \left\{ \left\langle \frac{\mathbf{v}^*}{\lambda}, \mathbf{v} \right\rangle - \mathbf{J}(\mathbf{v}) \right\} = \lambda \mathbf{J}^*\left(\frac{\mathbf{v}^*}{\lambda}\right). \end{aligned}$$

Property 3

Let $\mathbf{J} : \mathbf{V} \rightarrow \overline{\mathbb{R}}$ and $\mathbf{J}_\alpha(\mathbf{v}) = \mathbf{J}(\mathbf{v}) + \alpha$, where $\alpha \in \mathbb{R}$. Then

$$\mathbf{J}_\alpha^*(\mathbf{v}^*) = \mathbf{J}^*(\mathbf{v}^*) - \alpha.$$

Proof. It follows from the obvious relation

$$\sup_{\mathbf{v} \in \mathbf{V}} \{\langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{J}(\mathbf{v}) - \alpha\} = \sup_{\mathbf{v} \in \mathbf{V}} \{\langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{J}(\mathbf{v})\} - \alpha.$$

Property 4

Let $\mathbf{v}_0 \in \mathbf{V}$ and $\mathbf{G}(\mathbf{v}) = \mathbf{J}(\mathbf{v} - \mathbf{v}_0)$. Then

$$\mathbf{G}^*(\mathbf{v}^*) = \mathbf{J}^*(\mathbf{v}^*) + \langle \mathbf{v}^*, \mathbf{v}_0 \rangle.$$

Proof. Since

$$\begin{aligned} \sup_{\mathbf{v} \in \mathbf{V}} \{ \langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{J}(\mathbf{v} - \mathbf{v}_0) \} &= \sup_{\mathbf{w} \in \mathbf{V}} \{ \langle \mathbf{v}^*, \mathbf{w} + \mathbf{v}_0 \rangle - \mathbf{J}(\mathbf{w}) \} \\ &= \sup_{\mathbf{w} \in \mathbf{V}} \{ \langle \mathbf{v}^*, \mathbf{w} \rangle - \mathbf{J}(\mathbf{w}) \} + \langle \mathbf{v}^*, \mathbf{v}_0 \rangle = \mathbf{J}^*(\mathbf{v}^*) + \langle \mathbf{v}^*, \mathbf{v}_0 \rangle, \end{aligned}$$

we arrive at the required relation.

Property 5

If $\mathbf{G}(\mathbf{v}) = \min_{i=1,\dots,N} \{\mathbf{J}_i(\mathbf{v})\}$, then $\mathbf{G}^*(\mathbf{v}^*) = \max_{i=1,\dots,N} \{\mathbf{J}_i^*(\mathbf{v}^*)\}$.

Proof. We have

$$\begin{aligned}\mathbf{G}^*(\mathbf{v}^*) &= \sup_{\mathbf{v} \in \mathbf{V}} \{ \langle \mathbf{v}^*, \mathbf{v} \rangle - \min_{i=1,\dots,N} \{ \mathbf{J}_i(\mathbf{v}) \} \} \\ &= \sup_{\mathbf{v} \in \mathbf{V}} \{ \langle \mathbf{v}^*, \mathbf{v} \rangle + \max_{i=1,\dots,N} \{ -\mathbf{J}_i(\mathbf{v}) \} \} \\ &= \sup_{\mathbf{v} \in \mathbf{V}} \max_{i=1,\dots,N} \{ \langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{J}_i(\mathbf{v}) \} \\ &= \max_{i=1,\dots,N} \sup_{\mathbf{v} \in \mathbf{V}} \{ \langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{J}_i(\mathbf{v}) \} = \max_{i=1,\dots,N} \{ \mathbf{J}_i^*(\mathbf{v}^*) \}.\end{aligned}$$

Property 6

If $\mathbf{G}(\mathbf{v}) = \max_{i=1,\dots,N} \{\mathbf{J}_i(\mathbf{v})\}$, then $\mathbf{G}^*(\mathbf{v}^*) \leq \min_{i=1,\dots,N} \{\mathbf{J}_i^*(\mathbf{v}^*)\}$.

Proof. By definition, we have

$$\begin{aligned} \mathbf{G}^*(\mathbf{v}^*) &= \sup_{\mathbf{v} \in \mathbf{V}} \{ \langle \mathbf{v}^*, \mathbf{v} \rangle - \max_{i=1,\dots,N} \{ \mathbf{J}_i(\mathbf{v}) \} \} \\ &= \sup_{\mathbf{v} \in \mathbf{V}} \{ \langle \mathbf{v}^*, \mathbf{v} \rangle + \min_{i=1,\dots,N} \{ -\mathbf{J}_i(\mathbf{v}) \} \} \\ &= \sup_{\mathbf{v} \in \mathbf{V}} \min_{i=1,\dots,N} \{ \langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{J}_i(\mathbf{v}) \}. \end{aligned}$$

Now we apply $\sup \inf \leq \inf \sup$ relation to $\langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{J}_i(\mathbf{v})$. Then,

$$\mathbf{G}^*(\mathbf{v}^*) \leq \min_{i=1,\dots,N} \sup_{\mathbf{v} \in \mathbf{V}} \{ \langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{J}_i(\mathbf{v}) \} = \min_{i=1,\dots,N} \{ \mathbf{J}_i^*(\mathbf{v}^*) \}.$$

Subdifferential

Definition

The functional $\mathbf{J} \mathbf{V} \rightarrow \mathbb{R}$ is called subdifferentiable at \mathbf{v}_0 if there exists an affine minorant $\ell \in \mathbb{AM}(\mathbf{J})$ such that $\mathbf{J}(\mathbf{v}_0) = \ell(\mathbf{v}_0)$. A minorant with this property is called the **exact minorant** at \mathbf{v}_0 .

Obviously, any affine minorant exact for \mathbf{J} at \mathbf{v}_0 has the form

$$\ell(\mathbf{v}) = \langle \mathbf{v}^*, \mathbf{v} - \mathbf{v}_0 \rangle + \mathbf{J}(\mathbf{v}_0), \quad \ell(\mathbf{v}) \leq \mathbf{J}(\mathbf{v}), \quad \forall \mathbf{v} \in \mathbf{V}.$$

The element \mathbf{v}^* is called a **subgradient** of \mathbf{J} at \mathbf{v}_0 .

The set of all subgradients of \mathbf{J} at \mathbf{v}_0 forms a **subdifferential**, which is usually denoted by $\partial\mathbf{J}(\mathbf{v}_0)$. It may be empty or contain one element or infinitely many elements.

An important property of convex functionals follows directly from the above definition. For a convex functional \mathbf{J} at a point \mathbf{v}_0 where it is finite, the exact affine minorant is evidently **exist!**

In other words, there is at least one element $\mathbf{v}^* \in \partial\mathbf{J}(\mathbf{v}_0)$ that "creates" an affine minorant such that

$$\begin{aligned}\langle \mathbf{v}^*, \mathbf{v} \rangle - \alpha &\leq \mathbf{J}(\mathbf{v}), \quad \forall \mathbf{v} \in \mathbf{V}, \\ \langle \mathbf{v}^*, \mathbf{v}_0 \rangle - \alpha &= \mathbf{J}(\mathbf{v}_0).\end{aligned}$$

By subtracting, we obtain

$$\mathbf{J}(\mathbf{v}) - \mathbf{J}(\mathbf{v}_0) \geq \langle \mathbf{v}^*, \mathbf{v} - \mathbf{v}_0 \rangle.$$

The inequality (313) presents the **basic incremental relation for convex functionals**.

Compound functionals

Let \mathbf{J} and \mathbf{J}^* be a pair of mutually dual convex functionals.

The functional $\mathbf{D}_{\mathbf{J}} : \mathbf{V} \times \mathbf{V}^* \rightarrow \mathbb{R}$ of the form

$$\mathbf{D}_{\mathbf{J}}(\mathbf{v}, \mathbf{v}^*) := \mathbf{J}(\mathbf{v}) + \mathbf{J}^*(\mathbf{v}^*) - \langle \mathbf{v}^*, \mathbf{v} \rangle.$$

is called it the **compound functional** associated with these pair of functionals.

We will see that compound functionals play an important role in the a posteriori analysis of linear and nonlinear variational problems.

Compound functionals are always **nonnegative**. Indeed,

$$J^*(\mathbf{v}^*) = \sup_{\mathbf{v} \in \mathbf{V}} (\langle \mathbf{v}^*, \mathbf{v} \rangle - J(\mathbf{v})) \geq \langle \mathbf{v}^*, \mathbf{v} \rangle - J(\mathbf{v}) \quad \forall \mathbf{v} \in \mathbf{V}$$

and

$$J^*(\mathbf{v}^*) + J(\mathbf{v}) - \langle \mathbf{v}^*, \mathbf{v} \rangle \geq 0 \quad \forall \mathbf{v}, \mathbf{v}^*$$

Compound functionals may vanish only on special sets, where \mathbf{v} and \mathbf{v}^* satisfy certain relations.

Theorem

Let \mathbf{J} be a proper convex functional and \mathbf{J}^ be its polar. Then, the following two statements are equivalent:*

$$\mathbf{J}(\mathbf{v}) + \mathbf{J}^*(\mathbf{v}^*) - \langle \mathbf{v}^*, \mathbf{v} \rangle = 0, \quad (313)$$

$$\mathbf{v}^* \in \partial \mathbf{J}(\mathbf{v}) \text{ and } \mathbf{v} \in \partial \mathbf{J}^*(\mathbf{v}^*). \quad (314)$$

Relations (314) are also called **duality relations** for the pair $(\mathbf{v}, \mathbf{v}^*)$.

Proof.

Assume that $\mathbf{v}^* \in \partial \mathbf{J}(\mathbf{v})$, i.e.,

$$\mathbf{J}(\mathbf{w}) \geq \mathbf{J}(\mathbf{v}) + \langle \mathbf{v}^*, \mathbf{w} - \mathbf{v} \rangle, \quad \forall \mathbf{w} \in \mathbf{V}.$$

Hence,

$$\langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{J}(\mathbf{v}) \geq \langle \mathbf{v}^*, \mathbf{w} \rangle - \mathbf{J}(\mathbf{w}), \quad \forall \mathbf{w} \in \mathbf{V}$$

and, consequently,

$$\langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{J}(\mathbf{v}) \geq \sup_{\mathbf{w} \in \mathbf{V}} \{ \langle \mathbf{v}^*, \mathbf{w} \rangle - \mathbf{J}(\mathbf{w}) \} = \mathbf{J}^*(\mathbf{v}^*),$$

what leads to the conclusion that $\mathbf{J}^*(\mathbf{v}^*) + \mathbf{J}(\mathbf{w}) - \langle \mathbf{v}^*, \mathbf{w} \rangle \leq 0$.

But the left-hand side is nonnegative, so that we obtain

$$\mathbf{D}_{\mathbf{J}}(\mathbf{v}^*, \mathbf{v}) = 0.$$

Assume that $\mathbf{v} \in \partial \mathbf{J}^*(\mathbf{v}^*)$. Then

$$\mathbf{J}^*(\mathbf{w}^*) \geq \mathbf{J}^*(\mathbf{v}^*) + \langle \mathbf{w}^* - \mathbf{v}^*, \mathbf{v} \rangle,$$

and we continue similarly to the previous case:

$$\langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{J}^*(\mathbf{v}^*) \geq \langle \mathbf{w}^*, \mathbf{v} \rangle - \mathbf{J}^*(\mathbf{w}^*), \quad \forall \mathbf{w}^* \in \mathbf{V}^*,$$

$$\langle \mathbf{v}^*, \mathbf{v} \rangle - \mathbf{J}^*(\mathbf{v}^*) \geq \mathbf{J}^{**}(\mathbf{v}) = \mathbf{J}(\mathbf{v}).$$

Thus, we again arrive at the conclusion that it can only be if $\mathbf{D}_{\mathbf{J}}(\mathbf{v}^*, \mathbf{v}) = \mathbf{0}$.

Assume that $\mathbf{D}_J(\mathbf{v}^*, \mathbf{v}) = \mathbf{0}$. Since

$$\mathbf{J}^*(\mathbf{v}^*) = \sup_{\mathbf{w} \in \mathbf{V}} \{\langle \mathbf{v}^*, \mathbf{w} \rangle - \mathbf{J}(\mathbf{w})\},$$

we obtain

$$\mathbf{0} = \mathbf{J}(\mathbf{v}) + \mathbf{J}^*(\mathbf{v}^*) - \langle \mathbf{v}^*, \mathbf{v} \rangle \geq \mathbf{J}(\mathbf{v}) - \mathbf{J}(\mathbf{w}) - \langle \mathbf{v}^*, \mathbf{v} - \mathbf{w} \rangle, \quad \forall \mathbf{w} \in \mathbf{V}.$$

Rewrite this inequality in a more familiar form:

$$\mathbf{J}(\mathbf{w}) - \mathbf{J}(\mathbf{v}) \geq \langle \mathbf{v}^*, \mathbf{w} - \mathbf{v} \rangle, \quad \forall \mathbf{w} \in \mathbf{V},$$

which means that $\mathbf{J}(\mathbf{v}) + \langle \mathbf{v}^*, \mathbf{v} - \mathbf{w} \rangle$ is an exact affine minorant of \mathbf{J} (at \mathbf{v}) and, consequently, $\mathbf{v}^* \in \partial \mathbf{J}(\mathbf{v})$. The proof of the fact that $\mathbf{v}^* \in \partial \mathbf{J}^*(\mathbf{v}^*)$ is quite analogous.

Properties of compound functionals

First, we note that, $\mathbf{D}_G(\mathbf{y}, \mathbf{y}^*)$ is **convex** with respect to \mathbf{y} and \mathbf{y}^* , but, in general, $\mathbf{D}_G(\mathbf{y}, \mathbf{y}^*)$ is a **nonconvex** functional on $\mathbf{Y} \times \mathbf{Y}^*$.

This fact is easily observed in the simplest case $\mathbf{Y} = \mathbb{R}$ if set

$$G(y) = \frac{1}{\alpha}|y|^\alpha \quad G^*(y) = \frac{1}{\alpha^*}|y|^{\alpha^*}.$$

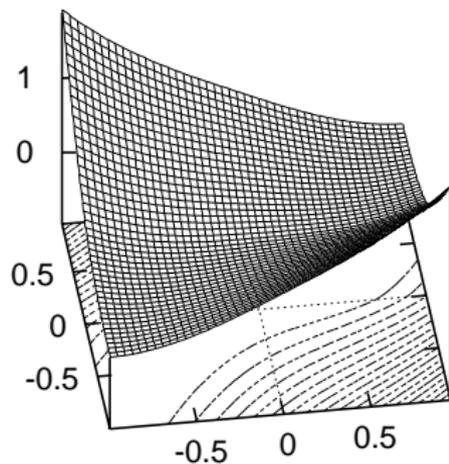
Only for $\alpha = 2$ we have a convex functional

$$\mathbf{D}_G(\mathbf{y}, \mathbf{y}^*) = \frac{1}{2}|y|^2 + \frac{1}{2}|y^*|^2 - yy^* = \frac{1}{2}(y - y^*)^2.$$

For other $\alpha \in (1, +\infty)$ \mathbf{D}_G is nonconvex on $\mathbb{R} \times \mathbb{R}$.

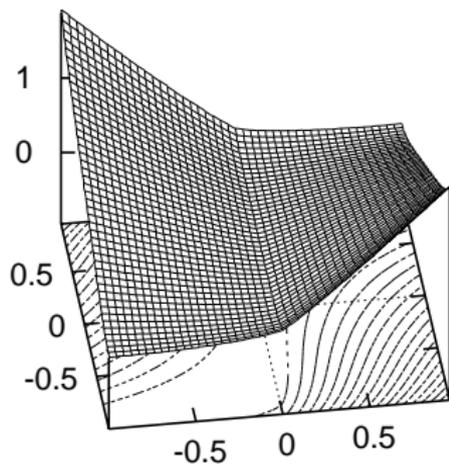
Example 1: $D(\xi_1, \xi_2) = \frac{1}{3}|\xi_1|^3 + \frac{2}{3}|\xi_2|^{3/2} - \xi_1\xi_2$

Compound functional on $\mathbb{R} \times \mathbb{R}$ and its level lines



Example 2: $D(\xi_1, \xi_2) = \frac{5}{6}|\xi_1|^{6/5} + \frac{1}{6}|\xi_2|^6 - \xi_1\xi_2$

Compound functional on $\mathbb{R} \times \mathbb{R}$ and its level lines



However, they have an important property, which is to some extent similar to convexity.

Theorem

For any $y_1, y_2 \in Y$ and $y_1^*, y_2^* \in Y^*$,

$$\mathbf{D}_G\left(\frac{y_1+y_2}{2}, \frac{y_1^*+y_2^*}{2}\right) \leq \frac{1}{4} \left(\mathbf{D}_G(y_1, y_1^*) + \mathbf{D}_G(y_1, y_2^*) + \mathbf{D}_G(y_2, y_1^*) + \mathbf{D}_G(y_2, y_2^*) \right)$$

Proof

From the definition it follows that

$$\begin{aligned} \mathbf{D}_G\left(\mathbf{y}, \frac{\mathbf{y}_1^* + \mathbf{y}_2^*}{2}\right) &= \mathbf{G}(\mathbf{y}) + \mathbf{G}^*\left(\frac{\mathbf{y}_1^* + \mathbf{y}_2^*}{2}\right) - \left\langle \frac{\mathbf{y}_1^* + \mathbf{y}_2^*}{2}, \mathbf{y} \right\rangle \\ &\leq \frac{1}{2} (\mathbf{D}_G(\mathbf{y}, \mathbf{y}_1^*) + \mathbf{D}_G(\mathbf{y}, \mathbf{y}_2^*)) \end{aligned}$$

and

$$\begin{aligned} \mathbf{D}_G\left(\frac{\mathbf{y}_1 + \mathbf{y}_2}{2}, \mathbf{y}^*\right) &= \mathbf{G}\left(\frac{\mathbf{y}_1 + \mathbf{y}_2}{2}\right) + \mathbf{G}^*(\mathbf{y}^*) - \left\langle \mathbf{y}^*, \frac{\mathbf{y}_1 + \mathbf{y}_2}{2} \right\rangle \\ &\leq \frac{1}{2} (\mathbf{D}_G(\mathbf{y}_1, \mathbf{y}^*) + \mathbf{D}_G(\mathbf{y}_2, \mathbf{y}^*)). \end{aligned}$$

Therefore,

$$\mathbf{D}_G\left(\frac{\mathbf{y}_1 + \mathbf{y}_2}{2}, \frac{\mathbf{y}_1^* + \mathbf{y}_2^*}{2}\right) \leq \frac{1}{2} \left(\mathbf{D}_G(\mathbf{y}_1, \frac{\mathbf{y}_1^* + \mathbf{y}_2^*}{2}) + \mathbf{D}_G(\mathbf{y}_2, \frac{\mathbf{y}_1^* + \mathbf{y}_2^*}{2}) \right).$$

and we arrive at the required estimate.

Important property

If \mathbf{G} and \mathbf{G}^* are Gateaux differentiable, then

$$\langle \mathbf{y}^* - \mathbf{G}'(\mathbf{y}), \mathbf{G}'(\mathbf{y}^*) - \mathbf{y} \rangle \geq \mathbf{D}_{\mathbf{G}}(\mathbf{y}, \mathbf{y}^*).$$

Note, that from this relation we conclude that $\mathbf{D}_{\mathbf{J}}$ vanishes if the duality relations are satisfied.

Uniformly convex functionals

Let a proper l.s.c. functional $\Upsilon : \mathbf{Y} \rightarrow \overline{\mathbb{R}}$ be subject to the conditions

$$\Upsilon(\mathbf{y}) \geq 0, \quad \forall \mathbf{y} \in \mathbf{Y}, \quad \Upsilon(\mathbf{y}) = 0 \iff \mathbf{y} = \mathbf{0}_{\mathbf{Y}}.$$

Definition

A convex functional $\mathbf{J} : \mathbf{Y} \rightarrow \overline{\mathbb{R}}$ is called **uniformly convex** in $\mathcal{B}(\mathbf{0}_{\mathbf{Y}}, \delta)$ if there exists a functional Υ_{δ} such that $\Upsilon_{\delta} \not\equiv 0$ and for all $\mathbf{y}_1, \mathbf{y}_2 \in \mathcal{B}(\mathbf{0}_{\mathbf{Y}}, \delta)$ the following inequality holds:

$$\mathbf{J}\left(\frac{\mathbf{y}_1 + \mathbf{y}_2}{2}\right) + \Upsilon_{\delta}(\mathbf{y}_1 - \mathbf{y}_2) \leq \frac{1}{2}(\mathbf{J}(\mathbf{y}_1) + \mathbf{J}(\mathbf{y}_2)). \quad (315)$$

The functional Υ_{δ} enforces standard convexity inequality. For this reason, it is called a **forcing** functional.

It is clear that any uniformly convex functional is convex in $\mathcal{B}(\mathbf{0}_Y, \delta)$. Now we establish two important inequalities that hold for uniformly convex functionals.

Theorem

If $\mathbf{J} : \mathbf{Y} \rightarrow \overline{\mathbb{R}}$ is uniformly convex in $\mathcal{B}(\mathbf{0}_Y, \delta)$ and Gâteaux differentiable in $\mathcal{B}(\mathbf{0}_Y, \delta)$, then for any $\mathbf{y}, \mathbf{z} \in \mathcal{B}(\mathbf{0}_Y, \delta)$ the following relations hold:

$$\mathbf{J}(\mathbf{z}) \geq \mathbf{J}(\mathbf{y}) + \langle \mathbf{J}'(\mathbf{y}), \mathbf{z} - \mathbf{y} \rangle + 2\Upsilon_\delta(\mathbf{z} - \mathbf{y})$$

and

$$\langle \mathbf{J}'(\mathbf{z}) - \mathbf{J}'(\mathbf{y}), \mathbf{z} - \mathbf{y} \rangle \geq 2\Upsilon_\delta(\mathbf{z} - \mathbf{y}) + 2\Upsilon_\delta(\mathbf{y} - \mathbf{z}).$$

Proof.

We have $\mathfrak{r}_\delta(\mathbf{z} - \mathbf{y}) \leq \frac{1}{2}\mathbf{J}(\mathbf{z}) + \frac{1}{2}\mathbf{J}(\mathbf{y}) - \mathbf{J}\left(\frac{\mathbf{z}+\mathbf{y}}{2}\right)$.

Since \mathbf{J} is convex and differentiable

$$\mathbf{J}\left(\frac{\mathbf{z}+\mathbf{y}}{2}\right) = \mathbf{J}\left(\mathbf{y} + \frac{\mathbf{z}-\mathbf{y}}{2}\right) \geq \mathbf{J}(\mathbf{y}) + \left\langle \mathbf{J}'(\mathbf{y}), \frac{\mathbf{z}-\mathbf{y}}{2} \right\rangle,$$

and, therefore,

$$2\mathfrak{r}_\delta(\mathbf{z} - \mathbf{y}) \leq \mathbf{J}(\mathbf{z}) - \mathbf{J}(\mathbf{y}) - \left\langle \mathbf{J}'(\mathbf{y}), \mathbf{z} - \mathbf{y} \right\rangle.$$

We can rewrite it replacing \mathbf{z} by \mathbf{y}

$$2\mathfrak{r}_\delta(\mathbf{y} - \mathbf{z}) \leq \mathbf{J}(\mathbf{y}) - \mathbf{J}(\mathbf{z}) + \left\langle \mathbf{J}'(\mathbf{z}), \mathbf{z} - \mathbf{y} \right\rangle$$

and obtain the second inequality. □

Deviations from the minimizer

Theorem

Let a functional \mathbf{J} be uniformly convex in $\mathcal{B}(\mathbf{0}_Y, \delta)$ and $\mathbf{y}_m \in \mathcal{B}(\mathbf{0}_Y, \delta)$ be the minimizer of \mathbf{J} .

$$\Upsilon_\delta(\mathbf{z} - \mathbf{y}_m) \leq \frac{1}{2} (\mathbf{J}(\mathbf{z}) - \mathbf{J}(\mathbf{y}_m)), \quad \forall \mathbf{z} \in \mathcal{B}(\mathbf{0}_Y, \delta). \quad (316)$$

Proof.

Since $\mathbf{J}\left(\frac{\mathbf{y}_m + \mathbf{z}}{2}\right) \geq \mathbf{J}(\mathbf{y}_m)$, we obtain

$$\begin{aligned} \Upsilon_\delta(\mathbf{z} - \mathbf{y}_m) &\leq \frac{1}{2} \mathbf{J}(\mathbf{y}_m) + \frac{1}{2} \mathbf{J}(\mathbf{z}) - \mathbf{J}\left(\frac{\mathbf{y}_m + \mathbf{z}}{2}\right) \leq \\ &\leq \frac{1}{2} (\mathbf{J}(\mathbf{z}) - \mathbf{J}(\mathbf{y}_m)). \end{aligned}$$



Estimate (316) is the first step in deriving a posteriori error estimates of the functional type by means of the variational techniques. It shows that deviations from the minimizer (measured in terms of the functional \mathfrak{T}_δ) are controlled by the difference of the functionals.

Corollary 1

Rewrite (315) in the form

$$\Upsilon_\delta(\mathbf{z} - \mathbf{y}_m) + \mathbf{J}\left(\frac{\mathbf{y}_m + \mathbf{z}}{2}\right) - \mathbf{J}(\mathbf{y}_m) \leq \frac{1}{2}(\mathbf{J}(\mathbf{z}) - \mathbf{J}(\mathbf{y}_m)).$$

By virtue of (316), we have

$$\mathbf{J}\left(\frac{\mathbf{y}_m + \mathbf{z}}{2}\right) - \mathbf{J}(\mathbf{y}_m) \geq 2\Upsilon_\delta\left(\frac{\mathbf{z} - \mathbf{y}_m}{2}\right)$$

and, therefore, we arrive at the strengthened estimate

$$\Upsilon_\delta(\mathbf{z} - \mathbf{y}_m) + 2\Upsilon_\delta\left(\frac{\mathbf{z} - \mathbf{y}_m}{2}\right) \leq \frac{1}{2}(\mathbf{J}(\mathbf{z}) - \mathbf{J}(\mathbf{y}_m)). \quad (317)$$

Corollary 2

Assume that \mathbf{J} is twice differentiable in the vicinity of \mathbf{y}_m and satisfies the finite increment relation

$$\begin{aligned} \mathbf{J}\left(\frac{\mathbf{y}_m + \mathbf{z}}{2}\right) = & \mathbf{J}(\mathbf{y}_m) + \left\langle \mathbf{J}'(\mathbf{y}_m), \frac{\mathbf{z} - \mathbf{y}_m}{2} \right\rangle + \\ & + \frac{1}{2} \left\langle \mathbf{J}''\left(\mathbf{y}_m + \xi \frac{\mathbf{z} + \mathbf{y}_m}{2}\right) \frac{\mathbf{z} - \mathbf{y}_m}{2}, \frac{\mathbf{z} - \mathbf{y}_m}{2} \right\rangle, \end{aligned}$$

where $\xi \in (0, 1)$. Since $\mathbf{J}'(\mathbf{y}_m) = \mathbf{0}_{\mathbf{y}^*}$, we have another estimate:

$$\begin{aligned} \mathfrak{r}_\delta(\mathbf{z} - \mathbf{y}_m) + \frac{1}{8} \left\langle \mathbf{J}''\left(\left(1 + \frac{\xi}{2}\right)\mathbf{y}_m + \frac{\xi}{2}\mathbf{z}\right) (\mathbf{z} - \mathbf{y}_m), \mathbf{z} - \mathbf{y}_m \right\rangle & \leq \\ & \leq \frac{1}{2} (\mathbf{J}(\mathbf{z}) - \mathbf{J}(\mathbf{y}_m)). \quad (318) \end{aligned}$$

Example 1

Consider a self-adjoint operator $\mathbf{A} \in \mathcal{L}(\mathbf{H}, \mathbf{H})$ defined on a Hilbert space \mathbf{H} with scalar product (\cdot, \cdot) . Assume that it satisfies the condition

$$\alpha_1 \|\mathbf{y}\|^2 \leq \mathbf{G}(\mathbf{y}) := (\mathbf{A}\mathbf{y}, \mathbf{y}) \leq \alpha_2 \|\mathbf{y}\|^2, \quad \forall \mathbf{y} \in \mathbf{H}.$$

For $\mathbf{J}(\mathbf{y}) = \mathbf{G}(\mathbf{y}) + (\ell, \mathbf{y})$, $\ell \in \mathbf{H}$ we have

$$\begin{aligned} \frac{1}{2}\mathbf{G}(\mathbf{y}) + \frac{1}{2}\mathbf{G}(\mathbf{z}) - \mathbf{G}\left(\frac{\mathbf{y} + \mathbf{z}}{2}\right) &= \\ &= \frac{1}{4}(\mathbf{A}\mathbf{y}, \mathbf{y}) + \frac{1}{4}(\mathbf{A}\mathbf{z}, \mathbf{z}) - \frac{1}{8}(\mathbf{A}(\mathbf{y} + \mathbf{z}), \mathbf{y} + \mathbf{z}) = \\ &= \frac{1}{8}(\mathbf{A}(\mathbf{z} - \mathbf{y}), (\mathbf{z} - \mathbf{y})), \end{aligned}$$

the functional \mathbf{G} is uniformly convex in any ball with

$$\Upsilon(\mathbf{z} - \mathbf{y}) = \frac{1}{8}(\mathbf{A}(\mathbf{z} - \mathbf{y}), (\mathbf{z} - \mathbf{y})).$$

Thus, from (316) we have

$$\frac{1}{8}(\mathbf{A}(\mathbf{z} - \mathbf{y}_m), (\mathbf{z} - \mathbf{y}_m)) \leq \frac{1}{2}(\mathbf{J}(\mathbf{z}) - \mathbf{J}(\mathbf{y}_m)), \quad \forall \mathbf{z}$$

However (318) gives a better estimate

$$\frac{1}{2}(\mathbf{A}(\mathbf{z} - \mathbf{y}_m), (\mathbf{z} - \mathbf{y}_m)) \leq \mathbf{J}(\mathbf{z}) - \mathbf{J}(\mathbf{y}_m). \quad (319)$$

Note that for quadratic type functionals this estimate holds as equality. Indeed,

$$\mathbf{J}(\mathbf{z}) - \mathbf{J}(\mathbf{y}_m) = (\mathbf{A}\mathbf{y}_m + \ell, \mathbf{z} - \mathbf{y}_m) + \frac{1}{2}(\mathbf{A}(\mathbf{z} - \mathbf{y}_m), \mathbf{z} - \mathbf{y}_m).$$

and the minimizer \mathbf{y}_m satisfies the relation

$$(\mathbf{A}\mathbf{y}_m + \ell, \mathbf{y}) = 0, \quad \forall \mathbf{y} \in \mathbf{Y}.$$

Therefore, (319) **holds as equality**.

Theorem

Let J_1 and J_2 be uniformly convex in $\mathcal{B}(0_Y, \delta)$ with functionals $\Upsilon_{1\delta}$ and $\Upsilon_{2\delta}$, respectively.

Then the functional

$$\mu_1 J_1 + \mu_2 J_2,$$

where $\mu_1, \mu_2 \geq 0$, is uniformly convex in $\mathcal{B}(0_Y, \delta)$ with

$$\Upsilon_\delta = \mu_1 \Upsilon_{1\delta} + \mu_2 \Upsilon_{2\delta}.$$

Proof.

The proposition follows directly from definition of uniform convexity . \square

Example 2

Consider the functional

$$\mathbf{J}(\mathbf{y}) = \frac{1}{2}(\mathbf{A}\mathbf{y}, \mathbf{y}) + (\ell, \mathbf{y}) + \Psi(\mathbf{y}),$$

where $\Psi(\mathbf{y})$ is a convex and l.s.c. functional. Applying the above Theorem with $\mu_1 = \mu_2 = 1$,

$$\mathbf{J}_1(\mathbf{y}) = \frac{1}{2}(\mathbf{A}\mathbf{y}, \mathbf{y}) + (\ell, \mathbf{y}) \quad \mathbf{J}_2(\mathbf{y}) = \Psi(\mathbf{y}),$$

we see that \mathbf{J} is uniformly convex with functional Υ defined in Example 1.

Theorem

Let \mathbf{J}_1 and \mathbf{J}_2 be uniformly convex in $\mathcal{B}(0_{\mathbf{Y}}, \delta)$ with functionals $\Upsilon_{1\delta}$ and $\Upsilon_{2\delta}$, respectively. Then the functional

$$\mathbf{J}(\mathbf{y}) = \max\{\mathbf{J}_1(\mathbf{y}), \mathbf{J}_2(\mathbf{y})\}$$

is uniformly convex in $\mathcal{B}(0_{\mathbf{Y}}, \delta)$ with

$$\Upsilon_{\delta} = \min\{\Upsilon_{1\delta}, \Upsilon_{2\delta}\}.$$

Proof. We have

$$\begin{aligned} \frac{1}{2}\mathbf{J}(\mathbf{y}) + \frac{1}{2}\mathbf{J}(\mathbf{z}) - \mathbf{J}\left(\frac{\mathbf{y}+\mathbf{z}}{2}\right) &= \frac{1}{2} \max\{\mathbf{J}_1(\mathbf{y}), \mathbf{J}_2(\mathbf{y})\} + \\ &+ \frac{1}{2} \max\{\mathbf{J}_1(\mathbf{z}), \mathbf{J}_2(\mathbf{z})\} - \max\left\{\mathbf{J}_1\left(\frac{\mathbf{y}+\mathbf{z}}{2}\right), \mathbf{J}_2\left(\frac{\mathbf{y}+\mathbf{z}}{2}\right)\right\}. \end{aligned}$$

Assume that

$$\max\left\{\mathbf{J}_1\left(\frac{\mathbf{y}+\mathbf{z}}{2}\right), \mathbf{J}_2\left(\frac{\mathbf{y}+\mathbf{z}}{2}\right)\right\} = \mathbf{J}_1\left(\frac{\mathbf{y}+\mathbf{z}}{2}\right).$$

Then

$$\begin{aligned} \frac{1}{2}(\mathbf{J}(\mathbf{y}) + \mathbf{J}(\mathbf{z})) - \mathbf{J}\left(\frac{\mathbf{y}+\mathbf{z}}{2}\right) &\geq \\ &\geq \frac{1}{2}(\mathbf{J}_1(\mathbf{y}) + \mathbf{J}_1(\mathbf{z})) - \mathbf{J}_1\left(\frac{\mathbf{y}+\mathbf{z}}{2}\right) \geq \mathfrak{r}_{1\delta}(\mathbf{z} - \mathbf{y}). \end{aligned}$$

If we have an opposite situation, i.e.,

$$\max \left\{ \mathbf{J}_1 \left(\frac{\mathbf{y}+\mathbf{z}}{2} \right), \mathbf{J}_2 \left(\frac{\mathbf{y}+\mathbf{z}}{2} \right) \right\} = \mathbf{J}_2 \left(\frac{\mathbf{y}+\mathbf{z}}{2} \right),$$

then

$$\frac{1}{2} \mathbf{J}(\mathbf{y}) + \frac{1}{2} \mathbf{J}(\mathbf{z}) - \mathbf{J} \left(\frac{\mathbf{y}+\mathbf{z}}{2} \right) \geq \mathfrak{T}_{2\delta}(\mathbf{z} - \mathbf{y}).$$

Thus, in both cases the lower bound is given by the functional

$$\min \left\{ \mathfrak{T}_{1\delta}(\mathbf{z} - \mathbf{y}), \mathfrak{T}_{2\delta}(\mathbf{z} - \mathbf{y}) \right\}.$$

Example 3. Power growth functionals

Let

$$\mathbf{G}(\mathbf{y}) = \frac{1}{\alpha} \int_{\Omega} |\mathbf{y}|^{\alpha} \, d\mathbf{x} \quad \mathbf{F}(\mathbf{v}) = \int_{\Omega} \mathbf{f} \mathbf{v} \, d\mathbf{x},$$

where $\alpha > 1$. Then Problem \mathcal{P} is to minimize the functional

$$\mathbf{J}_{\alpha}(\mathbf{v}) := \int_{\Omega} \left(\frac{1}{\alpha} |\nabla \mathbf{v}|^{\alpha} + \mathbf{f} \mathbf{v} \right) \, d\mathbf{x}$$

over the space $\mathbf{V} = \{\mathbf{v} \in \mathbf{H}^{\alpha}(\Omega) \mid \mathbf{v} = \mathbf{0} \text{ on } \partial\Omega\}$.

Problem \mathcal{P}^* is to maximize the functional

$$I_{\alpha^*}^*(\mathbf{y}^*) = -\frac{1}{\alpha^*} \int_{\Omega} |\mathbf{y}^*|^{\alpha^*} \, d\mathbf{x}$$

over the set

$$\mathbf{Q}_f^* = \left\{ \mathbf{y}^* \in \mathbf{Y}^* := \mathbf{L}^{\alpha^*}(\Omega, \mathbb{R}^n) \mid \int_{\Omega} \mathbf{y}^* \cdot \nabla \mathbf{w} \, d\mathbf{x} = \int_{\Omega} \mathbf{f} \mathbf{w} \, d\mathbf{x} \quad \forall \mathbf{w} \in \mathbf{V} \right\}.$$

For $\alpha \geq 2$ uniform convexity of $\mathbf{G}(\mathbf{y})$ follows from the first Clarkson's inequality

$$\int_{\Omega} \left| \frac{\mathbf{y}_1 + \mathbf{y}_2}{2} \right|^\alpha \mathbf{d}\mathbf{x} + \int_{\Omega} \left| \frac{\mathbf{y}_1 - \mathbf{y}_2}{2} \right|^\alpha \mathbf{d}\mathbf{x} \leq \frac{1}{2} \int_{\Omega} (|\mathbf{y}_1|^\alpha + |\mathbf{y}_2|^\alpha) \mathbf{d}\mathbf{x},$$

which is valid for all $\mathbf{y}_1, \mathbf{y}_2 \in \mathbf{Y}$.

See S. L. Sobolev. *Some Applications of Functional Analysis in Mathematical Physics*. Hence, we observe that in this case

$$\Upsilon_{\Theta}(\mathbf{z}) = \frac{1}{\alpha} \|\mathbf{z}\|_{\alpha, \Omega}^\alpha.$$

and

$$\frac{1}{\alpha} \int_{\Omega} |\nabla(\mathbf{v} - \mathbf{u})|^\alpha \mathbf{d}\mathbf{x} \leq \frac{1}{\alpha} (\mathbf{J}_{\Theta}(\mathbf{v}) - \mathbf{I}_{\Theta}^*(\mathbf{a}^*)), \quad \forall \mathbf{a}^* \in \mathbf{Q}_f^*.$$

For $1 < \alpha \leq 2$, the functional \mathbf{G} is also uniformly convex. This fact follows from the second Clarkson's inequality

$$\left(\int_{\Omega} \left(\frac{\mathbf{y}_1 + \mathbf{y}_2}{2} \right)^\alpha \mathbf{d}\mathbf{x} \right)^{\frac{1}{\alpha-1}} + \left(\int_{\Omega} \left(\frac{\mathbf{y}_1 - \mathbf{y}_2}{2} \right)^\alpha \mathbf{d}\mathbf{x} \right)^{\frac{1}{\alpha-1}} \leq \left(\frac{1}{2} \int_{\Omega} (|\mathbf{y}_1|^\alpha + |\mathbf{y}_2|^\alpha) \mathbf{d}\mathbf{x} \right)^{\frac{1}{\alpha-1}}.$$

However, in this case, the functional \mathfrak{T}_δ depends on the radius δ of a ball $\mathfrak{B}(\mathbf{0}_Y, \delta)$ that contains \mathbf{y}_1 and \mathbf{y}_2 , so that the estimate holds with

$$\mathfrak{T}_\delta(\mathbf{z}) = \delta^{\frac{\alpha-2}{\alpha-1}} \kappa \|\mathbf{z}\|_{\alpha, \Omega}^{\frac{\alpha}{\alpha-1}},$$

where $\kappa = \frac{1}{\kappa_0+1}$ and κ_0 is the integer part of $\frac{1}{\alpha-1}$.

Now we introduce a general scheme for deriving a posteriori error estimates by using duality theory of the calculus of variations. We consider variational problems of the form

$$\inf_{\mathbf{v} \in \mathbf{V}} \{ \mathbf{F}(\mathbf{v}) + \mathbf{G}(\mathbf{\Lambda v}) \},$$

where $\mathbf{F} : \mathbf{V} \rightarrow \mathbb{R}$ is a convex lower semicontinuous functional, $\mathbf{G} : \mathbf{Y} \rightarrow \mathbb{R}$ is a uniformly convex functional, \mathbf{V} and \mathbf{Y} are reflexive Banach spaces and $\mathbf{\Lambda} : \mathbf{V} \rightarrow \mathbf{Y}$ is a bounded linear operator.

General variational problem

Consider the general variational problem: find \mathbf{u} in a Banach space V such that

$$\mathbf{J}(\mathbf{u}, \mathbf{\Lambda u}) = \inf_{\mathbf{v} \in V} \mathbf{J}(\mathbf{v}, \mathbf{\Lambda v}), \quad (320)$$

where $\mathbf{J}(\mathbf{v}) = \mathbf{F}(\mathbf{v}) + \mathbf{G}(\mathbf{\Lambda v})$, \mathbf{F} is a convex, lower semicontinuous functional, \mathbf{G} is a uniformly convex functional and $\mathbf{\Lambda} : V \rightarrow Y$ is a bounded linear operator.

V and Y are reflexive Banach spaces endowed with the norms $\|\cdot\|_V$ and $\|\cdot\|$, respectively.

Dual spaces are denoted by \mathbf{V}^* and \mathbf{Y}^* with duality pairings $\langle \cdot, \cdot \rangle$ and $\langle\langle \cdot, \cdot \rangle\rangle$, respectively. The spaces \mathbf{Y} and \mathbf{Y}^* are endowed with the norms $\|\cdot\|$ and $\|\cdot\|_*$.

We assume that

$$\|\Lambda \mathbf{w}\| \geq \mathbf{c}_0 \|\mathbf{w}\|_{\mathbf{V}} \quad \forall \mathbf{w} \in \mathbf{V}, \quad (321)$$

where \mathbf{c}_0 is a positive constant independent of \mathbf{w} .

In addition to Λ , we introduce its conjugate $\Lambda^* : \mathbf{Y}^* \rightarrow \mathbf{V}^*$. This amounts to say that

$$\langle\langle \mathbf{y}^*, \Lambda \mathbf{v} \rangle\rangle = \langle \Lambda^* \mathbf{y}^*, \mathbf{v} \rangle \quad \forall \mathbf{y}^* \in \mathbf{Y}^*, \mathbf{v} \in \mathbf{V}. \quad (322)$$

$\mathbf{J}(\mathbf{v}, \Lambda \mathbf{v}) := \mathbf{F}(\mathbf{v}) + \mathbf{G}(\Lambda \mathbf{v})$ is assumed to be coercive on \mathbf{V} , i.e.

$$\mathbf{J}(\mathbf{v}, \Lambda \mathbf{v}) \rightarrow +\infty \quad \text{if} \quad \|\mathbf{v}\|_{\mathbf{V}} \rightarrow +\infty.$$

Primal and Dual Problems

Problem \mathcal{P} . Find $\mathbf{u} \in \mathbf{V}$ such that

$$\mathbf{J}(\mathbf{u}, \Lambda \mathbf{u}) = \inf \mathcal{P} := \inf_{\mathbf{v} \in \mathbf{V}} \mathbf{J}(\mathbf{v}, \Lambda \mathbf{v}). \quad (323)$$

The problem dual to (323 is (see e.g.

I. Ekeland and R. Temam *Convex analysis and variational problems*.
North-Holland, Amsterdam, 1976.)

Problem \mathcal{P}^* . Find $\mathbf{p}^* \in \mathbf{Y}^*$ such that

$$\begin{aligned} - \quad \mathbf{J}^*(\Lambda^* \mathbf{p}^*, -\mathbf{p}^*) &= \sup \mathcal{P}^* := \sup_{\mathbf{y}^* \in \mathbf{Y}^*} -\mathbf{J}^*(\Lambda^* \mathbf{y}^*, -\mathbf{y}^*), & (324) \\ \mathbf{J}^*(\Lambda^* \mathbf{y}^*, -\mathbf{y}^*) &:= \mathbf{F}^*(\Lambda^* \mathbf{y}^*) + \mathbf{G}^*(-\mathbf{y}^*), \end{aligned}$$

where \mathbf{F}^* and \mathbf{G}^* are the functionals conjugate of \mathbf{F} and \mathbf{G} , respectively.

Theorem (1)

If the functional \mathbf{F} is finite at some $\mathbf{u}_0 \in \mathbf{V}$ and the functional \mathbf{G} is continuous and finite at $\Lambda \mathbf{u}_0 \in \mathbf{Y}$, then there exists a minimizer \mathbf{u} to Problem \mathcal{P} and a maximizer \mathbf{p}^* to Problem \mathcal{P}^* . Besides,

$$\inf \mathcal{P} = \sup \mathcal{P}^* \quad (325)$$

and the following duality relations hold

$$\begin{aligned} \text{(i)} \quad & \mathbf{F}(\mathbf{u}) + \mathbf{F}^*(\Lambda^* \mathbf{p}^*) - \langle \Lambda^* \mathbf{p}^*, \mathbf{u} \rangle = 0, \\ \text{(ii)} \quad & \mathbf{G}(\Lambda \mathbf{u}) + \mathbf{G}^*(-\mathbf{p}^*) + \langle \mathbf{p}^*, \Lambda \mathbf{u} \rangle = 0. \end{aligned} \quad (326)$$

Above relations are equivalent to

$$\text{(i)} \quad \Lambda^* \mathbf{p}^* \in \partial \mathbf{F}(\mathbf{u}), \quad \text{(ii)} \quad -\mathbf{p}^* \in \partial \mathbf{G}(\Lambda \mathbf{u}).$$

Problems with uniformly convex functionals

We recall (see Lecture 4) that a continuous functional $\mathbf{G} : \mathbf{Y} \rightarrow \mathbb{R}$ is uniformly convex in a ball $\mathbf{B}(\mathbf{0}, \delta) := \{\mathbf{y} \in \mathbf{Y} \mid \|\mathbf{y}\| < \delta\}$ if there exists a continuous functional $\Phi_\delta : \mathbf{Y} \rightarrow \mathbb{R}_+$ such that $\Phi_\delta(\mathbf{y}) = \mathbf{0}$ only if $\mathbf{y} = \mathbf{0}_y$ is and

$$\mathbf{G}\left(\frac{\mathbf{y}_1 + \mathbf{y}_2}{2}\right) + \Phi_\delta(\mathbf{y}_2 - \mathbf{y}_1) \leq \frac{1}{2}(\mathbf{G}(\mathbf{y}_1) + \mathbf{G}(\mathbf{y}_2)) \quad \forall \mathbf{y}_1, \mathbf{y}_2 \in \mathbf{B}(\mathbf{0}, \delta).$$

Usually, Φ_δ is given by a continuous strictly increasing function of the norm $\|\mathbf{y}\|$.

General form of a posteriori estimates for uniformly convex variational problems was established in

S. Repin. A posteriori error estimation for variational problems with uniformly convex functionals, *Math. Comput.*, 69(230), 2000, 481-500.

General form of the functional a posteriori estimate

Theorem (2)

Assume that the above conditions on \mathbf{F} and \mathbf{G} are satisfied and

(i) \mathbf{G} is uniformly convex on a ball $B(0, \delta)$,

(ii) the solution \mathbf{u} of Problem \mathcal{P} and an element $\mathbf{v} \in \mathbf{V}$ are such, that

$\Lambda \mathbf{u}, \Lambda \mathbf{v} \in \mathbf{B}(0, \delta)$.

Then, for any $\mathbf{y}^* \in \mathbf{Y}^*$

$$\Phi_\delta(\Lambda(\mathbf{v} - \mathbf{u})) \leq \mathbf{M}_\oplus(\mathbf{v}, \mathbf{y}^*) := \mathbf{D}_\mathbf{F}(\Lambda^* \mathbf{y}^*, \mathbf{v}) + \mathbf{D}_\mathbf{G}(\mathbf{y}^*, \Lambda \mathbf{v}) \quad (327)$$

where

$$\mathbf{D}_\mathbf{F}(\Lambda^* \mathbf{y}^*, \mathbf{v}) := \frac{1}{2} (\mathbf{F}(\mathbf{v}) + \mathbf{F}^*(\Lambda^* \mathbf{y}^*) - \langle \Lambda^* \mathbf{y}^*, \mathbf{v} \rangle),$$

$$\mathbf{D}_\mathbf{G}(\mathbf{y}^*, \Lambda \mathbf{v}) := \frac{1}{2} (\mathbf{G}(\Lambda \mathbf{v}) + \mathbf{G}^*(-\mathbf{y}^*) + \langle \mathbf{y}^*, \Lambda \mathbf{v} \rangle).$$

Proof

Since \mathbf{F} is convex and \mathbf{G} is uniformly convex we obtain

$$\begin{aligned} \Phi_\delta(\Lambda(\mathbf{v} - \mathbf{u})) + \mathbf{G}(\Lambda(\frac{\mathbf{v} + \mathbf{u}}{2})) + \mathbf{F}(\frac{\mathbf{v} + \mathbf{u}}{2}) \leq \\ \frac{1}{2} \left[(\mathbf{F}(\mathbf{v}) + \mathbf{G}(\Lambda\mathbf{v})) + (\mathbf{F}(\mathbf{u}) + \mathbf{G}(\Lambda\mathbf{u})) \right]. \end{aligned}$$

The element \mathbf{u} is a minimizer, therefore

$$\mathbf{G}(\Lambda\mathbf{u}) + \mathbf{F}(\mathbf{u}) = \mathbf{J}(\mathbf{u}) \leq \mathbf{G}(\Lambda(\frac{\mathbf{u} + \mathbf{v}}{2})) + \mathbf{F}(\frac{\mathbf{u} + \mathbf{v}}{2})$$

and we have

$$\begin{aligned} \Phi_\delta(\Lambda(\mathbf{v} - \mathbf{u})) + \mathbf{G}(\Lambda\mathbf{u}) + \mathbf{F}(\mathbf{u}) \leq \\ \frac{1}{2} \left[(\mathbf{F}(\mathbf{v}) + \mathbf{G}(\Lambda\mathbf{v})) + (\mathbf{F}(\mathbf{u}) + \mathbf{G}(\Lambda\mathbf{u})) \right]. \end{aligned}$$

From the above we observe that

$$\begin{aligned}\Phi_\delta(\Lambda \mathbf{e}) &\leq \frac{1}{2} \left[(\mathbf{F}(\mathbf{v}) + \mathbf{G}(\Lambda \mathbf{v})) - (\mathbf{F}(\mathbf{u}) + \mathbf{G}(\Lambda \mathbf{u})) \right] = \\ &= \frac{1}{2} (\mathbf{J}(\mathbf{v}, \Lambda \mathbf{v}) - \mathbf{J}(\mathbf{u}, \Lambda \mathbf{u})) \quad \forall \mathbf{v} \in \mathbf{B}(\mathbf{0}, \delta).\end{aligned}$$

In view of Theorem 1,

$$\mathbf{J}(\mathbf{u}, \Lambda \mathbf{u}) = \inf \mathcal{P} = \sup \mathcal{P}^* = -\mathbf{F}^*(\Lambda^* \mathbf{p}^*) - \mathbf{G}^*(-\mathbf{p}^*).$$

Since \mathbf{p}^* is a solution of the dual problem, we have

$$-\mathbf{J}^*(\Lambda^* \mathbf{p}^*, -\mathbf{p}^*) \geq -\mathbf{J}^*(\Lambda^* \mathbf{y}^*, -\mathbf{y}^*) \quad \forall \mathbf{y}^* \in \mathbf{Y}^*,$$

so that

$$\mathbf{J}(\mathbf{u}, \Lambda \mathbf{u}) \geq -\mathbf{F}^*(\Lambda^* \mathbf{y}^*) - \mathbf{G}^*(-\mathbf{y}^*).$$

Therefore

$$\begin{aligned} \Phi_\delta(\Lambda \mathbf{e}) &\leq \frac{1}{2} (\mathbf{F}(\mathbf{v}) + \mathbf{G}(\Lambda \mathbf{v}) + \mathbf{F}^*(\Lambda^* \mathbf{p}^*) + \mathbf{G}^*(-\mathbf{p}^*)) \leq \\ &\leq \frac{1}{2} (\mathbf{F}(\mathbf{v}) + \mathbf{G}(\Lambda \mathbf{v}) + \mathbf{F}^*(\Lambda^* \mathbf{y}^*) + \mathbf{G}^*(-\mathbf{y}^*)). \end{aligned}$$

However, by (322) we observe that

$$\langle\langle \mathbf{y}^*, \Lambda \mathbf{v} \rangle\rangle - \langle \Lambda^* \mathbf{y}^*, \mathbf{v} \rangle = \mathbf{0} \quad \forall \mathbf{y}^* \in \mathbf{Y}^*, \mathbf{v} \in \mathbf{V}.$$

We add this zero term to the above relation and obtain the required estimate.

□

Comments

The right-hand side of (327) is the **sum of two compound functionals**

$$\mathbf{M}_F : \mathbf{V}^* \times \mathbf{V} \rightarrow \mathbb{R} \quad \text{and} \quad \mathbf{M}_G : \mathbf{Y}^* \times \mathbf{Y} \rightarrow \mathbb{R}.$$

They are nonnegative and vanishes if and only if \mathbf{v} and \mathbf{y}^* satisfy the relations (326)(i)–(ii).

Therefore, $\mathbf{M}_\oplus(\mathbf{v}, \mathbf{y}^*)$ is, in fact, a measure of the error in the **duality relations** for the pair $(\mathbf{v}, \mathbf{y}^*)$.

It vanishes if and only if $\mathbf{v} = \mathbf{u}$ and $\mathbf{y}^* = \mathbf{p}^*$.

Let the functional \mathbf{F} be uniformly convex on \mathbf{V} with a forcing functional φ_δ . Then the "forcing functional" has the form we have

$$\Phi_\delta(\Lambda \mathbf{e}) + \varphi_\delta(\mathbf{e}) \leq \frac{1}{2}(\mathbf{J}(\mathbf{v}, \Lambda \mathbf{v}) - \mathbf{J}(\mathbf{u}, \Lambda \mathbf{u})) \quad (328)$$

and, as a result, (327) is replaced by the strengthened estimate

$$\Phi_\delta(\Lambda \mathbf{e}) + \varphi_\delta(\mathbf{e}) \leq \mathbf{M}_\oplus(\mathbf{v}, \mathbf{y}^*) \quad \forall \mathbf{y}^* \in \mathbf{Y}^*. \quad (329)$$

It is not difficult to verify that

$$\begin{aligned}
 \mathbf{M}_{\oplus}(\mathbf{v}, \mathbf{y}^*) - \mathbf{M}_{\oplus}(\mathbf{v}, \mathbf{p}^*) &= \\
 &= \frac{1}{2} (\mathbf{F}(\mathbf{v}) + \mathbf{F}^*(\boldsymbol{\Lambda}^* \mathbf{y}^*) - \langle \boldsymbol{\Lambda}^* \mathbf{y}^*, \mathbf{v} \rangle + \mathbf{G}(\boldsymbol{\Lambda} \mathbf{v}) + \mathbf{G}^*(-\mathbf{y}^*) + \langle \langle \mathbf{y}^*, \boldsymbol{\Lambda} \mathbf{v} \rangle \rangle) - \\
 &\frac{1}{2} (\mathbf{F}(\mathbf{v}) + \mathbf{F}^*(\boldsymbol{\Lambda}^* \mathbf{p}^*) - \langle \boldsymbol{\Lambda}^* \mathbf{p}^*, \mathbf{v} \rangle + \mathbf{G}(\boldsymbol{\Lambda} \mathbf{v}) + \mathbf{G}^*(-\mathbf{p}^*) + \langle \langle \mathbf{p}^*, \boldsymbol{\Lambda} \mathbf{v} \rangle \rangle) = \\
 &= \mathbf{J}^*(\boldsymbol{\Lambda}^* \mathbf{y}^*, -\mathbf{y}^*) - \mathbf{J}^*(\boldsymbol{\Lambda}^* \mathbf{p}^*, -\mathbf{p}^*) \geq \mathbf{0}.
 \end{aligned}$$

Therefore, for any \mathbf{v} the right-hand side of (327) is minimal if $\mathbf{y}^* = \mathbf{p}^*$. Consequently, to make the estimate effective we have to find some \mathbf{y}^* close to \mathbf{p}^* in \mathbf{Y}^* . A simple way to obtain a function "close" to \mathbf{p}^* it to use duality relations. To this end, we set $\mathbf{y}^* = \boldsymbol{\sigma}^*(\mathbf{v})$, where

$$-\boldsymbol{\sigma}^*(\mathbf{v}) \in \partial \mathbf{G}(\boldsymbol{\Lambda} \mathbf{v}) \subset \mathbf{Y}^*.$$

In this case,

$$\mathbf{M}_G(\boldsymbol{\sigma}^*(\mathbf{v}), \boldsymbol{\Lambda}\mathbf{v}) = \mathbf{0}$$

and we get the estimate

$$\Phi_\delta(\boldsymbol{\Lambda}\mathbf{e}) \leq \mathbf{M}_F(\boldsymbol{\Lambda}^*\boldsymbol{\sigma}^*(\mathbf{v}), \mathbf{v}) \quad (330)$$

whose right-hand side depends on \mathbf{v} only.

However, the estimate (330) cannot be directly applied in one practically important case which we consider below.

Example. Diffusion problem with Robin conditions

Consider the variational problems for the functional

$$\mathbf{J}(\mathbf{v}, \nabla \mathbf{v}) = \int_{\Omega} \left(\frac{1}{2} |\nabla \mathbf{v}|^2 + \frac{\delta}{2} |\mathbf{v}|^2 \right) d\mathbf{x} + \int_{\partial_2 \Omega} \left(\frac{\alpha}{2} |\mathbf{v}|^2 - \mathbf{g}\mathbf{v} \right) d\mathbf{s}.$$

Our problem is to minimize \mathbf{J} on the set of functions vanishing at $\partial_1 \Omega$. Minimizer \mathbf{u} of this variational problem is related to the system

$$\begin{aligned} -\Delta \mathbf{u} + \delta \mathbf{u} &= \mathbf{0}, & \text{in } \Omega, \\ \frac{\partial \mathbf{u}}{\partial \mathbf{n}} + \alpha \mathbf{u} - \mathbf{g} &= \mathbf{0}, & \text{on } \partial_2 \Omega. \end{aligned}$$

On $\partial_2 \Omega$ the solution satisfies the so-called **Robin** boundary condition. Let us show that the respective functional a posteriori estimate for the problem with Robin type boundary conditions easily follows from the above general estimate.

We set

$$\Lambda \mathbf{v} := \nabla \mathbf{v},$$
$$\mathbf{G}(\Lambda \mathbf{w}) = \int_{\Omega} \frac{1}{2} |\nabla \mathbf{v}|^2 \mathbf{d}\mathbf{x}$$

and

$$\mathbf{F}(\mathbf{v}) = \int_{\Omega} \frac{\delta}{2} |\mathbf{v}|^2 \mathbf{d}\mathbf{x} + \int_{\partial_2 \Omega} \left(\frac{\alpha}{2} |\mathbf{v}|^2 - \mathbf{g}\mathbf{v} \right) \mathbf{d}\mathbf{s}.$$

Since

$$\int_{\Omega} \mathbf{y}^* \cdot \nabla \mathbf{v} \, dx = \int_{\Omega} -\operatorname{div} \mathbf{y}^* \mathbf{v} \, dx + \int_{\partial_2 \Omega} (\mathbf{y}^* \cdot \mathbf{n}) \mathbf{v} \, ds,$$

we observe that $\mathbf{\Lambda}^* \mathbf{y}^* = \{-\operatorname{div} \mathbf{y}^* \mid_{\Omega}, \mathbf{y}^* \cdot \mathbf{n} \mid_{\partial_2 \Omega}\}$.

In the considered case,

$$\langle \mathbf{y}^*, \mathbf{y} \rangle := \int_{\Omega} \mathbf{y}^* \cdot \mathbf{y} \, dx;$$

$$\mathbf{G}^*(-\mathbf{y}^*) = \sup_{\mathbf{y}} \int_{\Omega} (-\mathbf{y}^* \cdot \mathbf{y} - \frac{1}{2} |\mathbf{y}|^2) \, dx = \int_{\Omega} \frac{1}{2} |\mathbf{y}^*|^2 \, dx.$$

Therefore,

$$\mathbf{G}(\mathbf{\Lambda}\mathbf{v}) + \mathbf{G}^*(-\mathbf{y}^*) + \langle \mathbf{y}^*, \mathbf{\Lambda}\mathbf{v} \rangle = \int_{\Omega} \left(\frac{1}{2} |\nabla \mathbf{v}|^2 + \frac{1}{2} |\mathbf{y}^*|^2 + \nabla \mathbf{v} \cdot \mathbf{y}^* \right) \mathbf{d}\mathbf{x}.$$

Next, in general,

$$\langle \langle \mathbf{\Lambda}^* \mathbf{y}^*, \mathbf{v} \rangle \rangle = \langle -\mathbf{div} \mathbf{y}^*, \mathbf{v} \rangle_{\mathbf{H}^{-1}(\Omega)} + \langle \mathbf{y}^* \cdot \mathbf{n}, \mathbf{v} \rangle_{\mathbf{H}^{-1/2}(\partial_2 \Omega)}.$$

However, if we assume that \mathbf{y}^* is sufficiently regular, then

$$\langle \langle \mathbf{\Lambda}^* \mathbf{y}^*, \mathbf{v} \rangle \rangle = \int_{\Omega} -\mathbf{div} \mathbf{y}^* \mathbf{v} \mathbf{d}\mathbf{x} + \int_{\partial_2 \Omega} \mathbf{y}^* \cdot \mathbf{n} \mathbf{v} \mathbf{d}\mathbf{s}.$$

Now,

$$\begin{aligned}
 \mathbf{F}^*(\boldsymbol{\Lambda}^* \mathbf{y}^*) &= \sup_{\mathbf{v}} \left\{ \int_{\Omega} -\operatorname{div} \mathbf{y}^* \mathbf{v} \, dx + \int_{\partial_2 \Omega} \mathbf{y}^* \cdot \mathbf{n} \, v \, ds - \mathbf{F}(\mathbf{v}) \right\} = \\
 &\sup_{\mathbf{v}} \left\{ \int_{\Omega} -\operatorname{div} \mathbf{y}^* \mathbf{v} \, dx + \int_{\partial_2 \Omega} \mathbf{y}^* \cdot \mathbf{n} \, v \, ds - \int_{\Omega} \frac{\delta}{2} |\mathbf{v}|^2 \, dx - \int_{\partial_2 \Omega} \left(\frac{\alpha}{2} |\mathbf{v}|^2 - \mathbf{g} \mathbf{v} \right) \, ds \leq \right. \\
 &\left. \sup_{\mathbf{v} \in \mathbf{L}_2(\Omega)} \int_{\Omega} \left(-\operatorname{div} \mathbf{y}^* \mathbf{v} - \frac{\delta}{2} |\mathbf{v}|^2 \right) \, dx + \sup_{\varrho \in \mathbf{L}_2(\partial_2 \Omega)} \int_{\partial_2 \Omega} \left((\mathbf{y}^* \cdot \mathbf{n}) \varrho - \frac{\alpha}{2} |\varrho|^2 + \mathbf{g} \varrho \right) \, ds \right.
 \end{aligned}$$

$$\sup_{\mathbf{v} \in L_2(\Omega)} \int_{\Omega} (-\operatorname{div} \mathbf{y}^* \cdot \mathbf{v} - \frac{\delta}{2} |\mathbf{v}|^2) \mathbf{d}\mathbf{x} = \int_{\Omega} \frac{1}{2\delta} |\operatorname{div} \mathbf{y}^*|^2 \mathbf{d}\mathbf{x},$$

$$\sup_{\varrho \in L_2(\partial_2 \Omega)} \int_{\partial_2 \Omega} ((\mathbf{y}^* \cdot \mathbf{n}) \varrho - \frac{\alpha}{2} |\varrho|^2 + \mathbf{g} \varrho) \mathbf{d}\mathbf{s} = \int_{\partial_2 \Omega} \frac{1}{2\alpha} |\mathbf{y}^* \cdot \mathbf{n} + \mathbf{g}|^2 \mathbf{d}\mathbf{s}.$$

Hence,

$$\mathbf{F}^*(\Lambda^* \mathbf{y}^*) \leq \int_{\Omega} \frac{1}{2\delta} |\operatorname{div} \mathbf{y}^*|^2 \mathbf{d}\mathbf{x} + \int_{\partial_2 \Omega} \frac{1}{2\alpha} |\mathbf{y}^* \cdot \mathbf{n} + \mathbf{g}|^2 \mathbf{d}\mathbf{s}.$$

Now,

$$\begin{aligned} \mathbf{F}(\mathbf{v}) &= \int_{\Omega} \frac{\delta}{2} |\mathbf{v}|^2 \mathbf{d}\mathbf{x} + \int_{\partial_2 \Omega} \left(\frac{\alpha}{2} |\mathbf{v}|^2 - \mathbf{g}\mathbf{v} \right) \mathbf{d}\mathbf{s}, \\ \langle\langle \boldsymbol{\Lambda}^* \mathbf{y}^*, \mathbf{v} \rangle\rangle &= \int_{\Omega} -\operatorname{div} \mathbf{y}^* \mathbf{v} \mathbf{d}\mathbf{x} + \int_{\partial_2 \Omega} (\mathbf{y}^* \cdot \mathbf{n}) \mathbf{v} \mathbf{d}\mathbf{s}, \\ \mathbf{F}^*(\boldsymbol{\Lambda}^* \mathbf{y}^*) &\leq \int_{\Omega} \frac{1}{2\delta} |\operatorname{div} \mathbf{y}^*|^2 \mathbf{d}\mathbf{x} + \int_{\partial_2 \Omega} \frac{1}{2\alpha} |\mathbf{y}^* \cdot \mathbf{n} + \mathbf{g}|^2 \mathbf{d}\mathbf{s}. \end{aligned}$$

Therefore,

$$\begin{aligned} \mathbf{F}(\mathbf{v}) + \mathbf{F}^*(\boldsymbol{\Lambda}^* \mathbf{y}^*) - \langle\langle \boldsymbol{\Lambda}^* \mathbf{y}^*, \mathbf{v} \rangle\rangle &\leq \int_{\Omega} \frac{1}{2\delta} (\operatorname{div} \mathbf{y}^* + \delta \mathbf{v})^2 \mathbf{d}\mathbf{x} + \\ &\int_{\partial_2 \Omega} \left(\frac{\alpha}{2} |\mathbf{v}|^2 + \frac{1}{2\alpha} |\mathbf{y}^* \cdot \mathbf{n} + \mathbf{g}|^2 - (\mathbf{y}^* \cdot \mathbf{n} + \mathbf{g}) \mathbf{v} \right) \mathbf{d}\mathbf{s}. \end{aligned}$$

We obtain

$$\mathbf{F}(\mathbf{v}) + \mathbf{F}^*(\boldsymbol{\Lambda}^* \mathbf{y}^*) - \langle \langle \boldsymbol{\Lambda}^* \mathbf{y}^*, \mathbf{v} \rangle \rangle \leq \int_{\Omega} \frac{1}{2\delta} (\mathbf{div} \mathbf{y}^* + \delta \mathbf{v})^2 \mathbf{d}\mathbf{x} + \int_{\partial_2 \Omega} \frac{1}{2\alpha} |\mathbf{y}^* \cdot \mathbf{n} + \mathbf{g} - \alpha \mathbf{v}|^2 \mathbf{d}\mathbf{s}.$$

$$\mathbf{G}(\boldsymbol{\Lambda} \mathbf{v}) + \mathbf{G}^*(-\mathbf{y}^*) + \langle \mathbf{y}^*, \boldsymbol{\Lambda} \mathbf{v} \rangle = \int_{\Omega} \frac{1}{2} |\nabla \mathbf{v} + \mathbf{y}^*|^2 \mathbf{d}\mathbf{x}.$$

Two terms above give the error Majorant.

We observe that the Majorant vanishes if and only if

$$\begin{aligned}\mathbf{div} \mathbf{y}^* + \delta \mathbf{v} &= \mathbf{0} && \text{in } \Omega, \\ \mathbf{y}^* \cdot \mathbf{n} + \mathbf{g} - \alpha \mathbf{v} &= \mathbf{0} && \text{on } \partial_2 \Omega, \\ \mathbf{y}^* &= -\nabla \mathbf{v} && \text{in } \Omega.\end{aligned}$$

These relations mean that

$$\begin{aligned}-\Delta \mathbf{v} + \delta \mathbf{v} &= \mathbf{0} && \text{in } \Omega, \\ \frac{\partial \mathbf{v}}{\partial \mathbf{n}} + \alpha \mathbf{v} &= \mathbf{g} && \text{on } \partial_2 \Omega,\end{aligned}$$

i.e., since \mathbf{v} vanishes at $\partial_1 \Omega$ it is but the exact solution.



M. Ainsworth and J. T. Oden. A posteriori error estimation in the finite element method, *Numer. Math.*, 60(1992) 429-463.



M. Ainsworth and J. T. Oden. A unified approach to a posteriori error estimation using element residual methods, *Numer. Math.*, 65(1993) 23-50.



M. Ainsworth and J. T. Oden, *A posteriori error estimation in finite element analysis*, Wiley and Sons, New York, 2000.



M. Ainsworth, J. T. Oden and C. Y. Lee. Local a posteriori error estimators for variational inequalities, *Numer. Methods for PDE*, 9(1993), 23-33.



M. Ainsworth, J. Z. Zhu, A. W. Craig and O. C. Zienkiewicz. Analysis of the Zienkiewicz-Zhu a posteriori error estimator in the finite element method, *Int. J. Numer. Methods Engrg.*, 28(1989), 2161-2174.



M. Amara, M. Ben Younes and C. Bernardi. Error indicators for the Navier-Stokes equations in stream function and vorticity formulation, *Numer. Math.*, 80(1998), 181-206.

-  A. Arkhipova. On the best possible smoothness of the problem with two-side constraints, *Vestnik Leningr. Univ., ser. Math.*, 7 (1984), 5-9 (in Russian).
-  A. M. Arthurth. Complementary variational principles, Clarendon Press, Oxford, 1980.
-  G. Astarita and G. Marrucci. *Principles of non Newtonian fluid mechanics*. McGraw Hill, London, 1974.
-  J.P.Aubin. *Approximation of elliptic boundary value problems*. Wiley-Interscience, New York-London-Sydney, 1972.
-  G. Auchmuty. A posteriori error estimates for linear equations, *Numer. Math.*, 61(1992), 1-6.
-  O. Axelsson. *Iterative solution methods*. Cambridge University Press, Cambridge, 1994.
-  I. Babuška. Courant element: before and after, in *Fifty years of Courant element*, M. Křížek, P. Neittaanmäki and R. Stenberg (Eds.), Marcel Dekker, 1994, 37-51.



I. Babuška. The finite element method with Lagrangian multipliers, *Numer. Math.*, 20(1973), 179-192.



I. Babuška, R. Duran and R. Rodriguez. Analysis of the efficiency of a posteriori error estimator for linear triangular elements, *SIAM J. Numer. Anal.* 29(1992), 947-964.



I. Babuška, F. Ihlenburg, T. Strouboulis and S. K. Gangaraj. A posteriori error estimation for finite element solutions on Helmholtz' equation—Part II: estimation of the pollution error, *Int. J. Numer. Meth. Engrg.*, 40(1997), 3883-3900.



I. Babuška and J. E. Osborn. Can a finite element method perform arbitrarily badly?, *Math. Comput.*, 69(2000), 230, 443-462.



I. Babuška and W. C. Rheinboldt. A-posteriori error estimates for the finite element method. *Internat. J. Numer. Meth. Engrg.*, 12(1978) 1597-1615.



I. Babuška and W. C. Rheinboldt. Error estimates for adaptive finite element computations. *SIAM J. Numer. Anal.*, 15(1978), 736-754.



I. Babuška and R. Rodriguez. The problem of the selection of an a posteriori error indicator based on smoothing techniques, *Internat. J. Numer. Meth. Engrg.*, 36(1993), 539-567.



I. Babuška and T. Strouboulis. *The finite element method and its reliability*. The Clarendon Press, Oxford University Press, New York, 2001.



I. Babuška, T. Strouboulis, C. S. Upadhyay and S. K. Gangaraj. A posteriori estimation and adaptive control of the pollution-error in the h -version of the FEM, *Int. J. Numer. Meth. Engrg.*, 38(1995), 4207-4235.



I. Babuška, T. Strouboulis, S. K. Gangaraj and C. S. Upadhyay. Pollution-error in the h -version of the FEM and the local quality of recovered derivatives, *Comput. Methods Appl. Mech. Engrg.*, 140(1997), 1-37.



I. Babuška, T. Strouboulis, A. Mathur and C. S. Upadhyay. Pollution error in the h -version of the FEM and the local quality of a posteriori error estimators, *Finite. Elem. Anal. and Des.*, 17(1994), 273-321.



I. Babuška, T. Strouboulis and C. S. Upadhyay. A model study of the quality of a posteriori error estimators for linear elliptic problems. Error estimation in the interior of patchwise uniform grids of triangles, *Comp. Meth. Appl. Mech. Engrg.* , 114(1994), 307-378.



C. Baiocchi and A. Capelo. *Variational and quasivariational inequalities. Applications to free boundary problems.* Wiley and Sons, New York, 1984.



J. M. Ball and B. D. James. Fine phase mixtures as minimizers of energy, *Arch. Rational Mech. Anal.*, 100(1987), 13-52.



P. Baumann and D. Phillips. A nonconvex variational problem related to change of phase, *Appl. Math. Optim.*, 21(1990), 113-138.



R. E. Bank and R. K. Smith. Mesh smoothing using a posteriori error estimates, *SIAM J. Numer. Anal.*, 34(1997), 3, 979-997.



R. E. Bank and A. Weiser. Some a posteriori error estimators for elliptic partial differential equations, *Math. Comput.*, 44(1985), 283-301.



R. E. Bank and B. D. Welfert. A posteriori error estimates for the Stokes problem, *SIAM J. Numer. Anal.*, 28(1991), 591-623.



W. Bangerth and R. Rannacher. *Adaptive finite element methods for differential equations*. Birkhäuser, Berlin, 2003.



R. Becker and R. Rannacher. A feed-back approach to error control in finite element methods: Basic approach and examples, *East-West J. Numer. Math.*, 4(1996), 237-264.



R. Becker, H. Kapp and R. Rannacher. Adaptive finite element methods for optimal control of partial differential equations: Basic concept. *SIAM J. Control Optim.*, 39(2000),1, 113-132.



A. Bermúdez, R. Durán and R. Rodríguez. Finite element analysis of compressible and incompressible fluid-solid systems, *Math. Comput.*, 67(1998), 221, 111-136.



B. Boroomand and O. C. Zienkiewicz, Recovery by equilibrium in patches (REP), *Int. J. Numer. Meth. Engrg.*, 40(1997), 137-164.

-  B. Boroomand and O. C. Zienkiewicz, An improved REP recovery and the effectivity robustness test, *Int. J. Numer. Meth. Engrg.*, 40(1997), 3247-3277.
-  S. W. Brady and A. R. Elcrat. Some results on a posteriori error estimation for approximate solution of second order elliptic problems, *SIAM J. Numer. Anal.*, 16(1979), 6, 877-889.
-  D. Braess. *Finite elements*. Cambridge University Press, Cambridge, 1997.
-  J.H. Bramble, R.D. Lazarov and J. E. Pasciak. Least-squares for second-order elliptic problems. *Comput. Methods Appl. Mech. Engrg.* , 152(1998), 195-210.
-  J. H. Bramble and R. S. Falk. A mixed-Lagrange multiplier finite element method for the polyharmonic equation, *RAIRO Model. Math. Anal. Numeric.* 19(4), 1985, 519-557.
-  H. Brezis. Problèmes unilatéraux, *J. Math. Pures Appl.*, 9(1971, 1, 1-168.



Brezis H., Kinderlehrer D. The smoothness of solutions to nonlinear variational inequalities, *Indiana Univ. Math J.*, 23(1974), 831-844.



H. Brezis and M. Sibony. Equivalence de deux inéquations variationnelles et applications, *Arch. Rat. Mech. Anal.*, 41(1971), 254-265.



F. Brezzi, On the existence, uniqueness and approximation of saddle-point problems arising from Lagrange multipliers, *R.A.I.R.O., Annal. Numer.*, 8 (1974), R-2, 129-151.



F. Brezzi, J. J. Douglas, R. Duran and M. Fortin. Mixed finite elements for second order elliptic problems in three variables, *Numer. Math.*, 51(1987), 2, 237-250.



F. Brezzi and M. Fortin. *Mixed and hybrid finite element methods*, Springer Series in Computational Mathematics, 15, New York, 1991.



F. Brezzi, M. Fortin and R. Stenberg. Error analysis of mixed-interpolated elements for Reissner-Midlin plates, *Math. Models and Methods in Applied Sciences*, 1(1991), 125-151.



F. Brezzi, C. Johnson, B. Mercier. Analysis of a mixed finite element method for elasto-plastic plates, *Mathematics of Computation*, (31)1977, 140, 809-817. Brezzi, F.; Pitkaranta, J. On the stabilization of finite element approximations of the Stokes equations. Efficient solutions of elliptic systems (Kiel, 1984), 11–19, *Notes Numer. Fluid Mech.*, 10, Vieweg, Braunschweig, 1984.



J. Bonvin, M. Picasso, and R. Stenberg, *GLS and EVSS methods for a three-field Stokes problem arising from viscoelastic flows*, *Comput. Methods Appl. Mech. Engrg.* 190(2001), 3893–3914.



I. G. Bubnov. *Selected Works*. Sudpromgiz, Leningrad (1956) (in Russian).



H. Buss and S. Repin. A posteriori error estimates for boundary-value problems with obstacles. In *Proceedings of 3d European Conference on Numerical Mathematics and Advanced Applications, Jyväskylä, 1999*, 162-170, World Scientific, 2000.



G. F. Carey and D. L. Humphrey. Mesh refinement and iterative solution methods for finite element computations, *Int. J. Numer. Meth. Engrg.*, 17(1981), 1717-1734.



G. F. Carey and A. I. Pehlivanov. Local error estimation and adaptive remeshing scheme for least-squares mixed finite elements, *Comput. Methods Appl. Mech. Engrg.*, 150(1997), 125-131.



C. Carstensen. A posteriori estimate for the mixed finite element method, *Math. Comput.*, 66(1997), 218, 465-476.



C. Carstensen. Quasi-interpolation and a posteriori error analysis of finite element methods, *Mathematical Modelling in Numerical Analysis*, 33(1999), 6, 1187-1202.



Carstensen, C.; Bartels, S. Each averaging technique yields reliable a posteriori error control in FEM on unstructured grids. I: Low order conforming, nonconforming, and mixed FEM, *Math. Comp.*, 71(2002), 239, 945-969.

-  C. Carstensen and S. A. Funken. Fully reliable localized error control in the FEM, *SIAM J. Sci. Comput.*, 21(2000), 4, 1465-1484.
-  C. Carstensen and S. A. Funken. Constants in Clément-interpolation error and residual based a posteriori error estimates in Finite Element Methods, *East-West J. Numer. Math.*, 8(2000), 3, 153-175.
-  C. Carstensen and R. Verfürth. Edge residuals dominate a posteriori error estimates for low order finite element methods, *SIAM J. Numer. Anal.*, 36 (1999), 5, 1571-1587.
-  A. Charbonneau, K. Dossou and R. Pierre. A residual-based a posteriori error estimator for Ciarlet–Raviart formulation of the first biharmonic problem, *Numer. Meth. for PDE's*, 13(1997), 93-111.
-  Z. Chen and R. H. Nochetto. Residual type a posteriori error estimates for elliptic obstacle problems, *Numer. Math.*, 84(2000), 527-548.
-  X. Cheng, W. Han and H. Huang. Analysis of some mixed elements for the Stokes problem, *Journal of Computational and Applied Mathematics*, 85(1997), 19-35.

-  M. Chipot. Regularity for the two obstacles problem. In *Free boundary problems, Vol. II (Pavia, 1979)*, Ist. Naz. Alta Mat. Francesco Severi, Rome, 1980, 135-140.
-  S. S. Chow. Finite element error estimates for non-linear elliptic equations of monotone type, *Numer. Math.*, 54(1989), 373-393.
-  S.-S. Chow, G. F. Carey and R. D. Lazarov. Natural and post-processed superconvergence in semilinear problems, *Numer. Meth. PDE*, 7(1991), 245-259.
-  P. G. Ciarlet. *The finite element method for elliptic problems*. North Holland, New York, Oxford, 1978.
-  Ph. Clément. Approximations by finite element functions using local regularization, *RAIRO Anal. Numér.*, 9(1975), R-2, 77-84.
-  *Funktionanalysis und numerische mathematik*. Springer-Verlag, Berlin, 1964.
-  P. Coorevits, J.-P. Dumeau and J.-P. Pelle. Error estimator and adaptivity for three-dimensional finite element analysis. In *Advances in*

Adaptive Computational Methods in Mechanics, Ed. P. Ladev eze and J. T. Oden, 443-457, Elsevier, New York, 1998.



R. Courant. Variational methods for some problems of equilibrium and vibration, *Bulletin of AMS*, 49(1943), 1-23.



Computational Methods in the mechanics of Fracture, Ed. S.N. Atluri, North-Holland, New York, 1986.



P. Destuynder and B. M etivet. Explicit error bounds of a conforming finite element method, *Math. Comput.*, 68(1999), 1379-1396.



W. D orfler and M. Rumpf. An adaptive strategy for elliptic problems including a posteriori controlled boundary approximation, *Math. Comput.*, 67(1998), 224, 1361-1362.



J. Douglas, Jr. T. Dupont and L. B. Wahlbin. The stability in L^q of the L^2 -projection into finite element function spaces, *Numer. Math.*, 23(1975), 193-197.



R. Duran, M. A. Muschietti and R. Rodriguez. On the asymptotic exactness of error estimators for linear triangle elements, *Numer. Math.*, 59(1991), 107-127.



R. Duran and R. Rodriguez. On the asymptotic exactness of Bank–Weiser’s estimator, *Numer. Math.*, 62(1992), 297-303.



G. Duvant and J.-L. Lions. *Les inequations en mecanique et en physique*, Dunod, Paris, 1972.



D. C. Drucker and W. Prager. Soil mechanics and plastic analysis or limit design. *Quart. Appl. Math.*, 10(1952), 157-165.



I. Ekeland and R. Temam. *Convex analysis and variational problems*. North-Holland, Amsterdam, 1976.



H.W. Engl and O. Scherzer. Convergence rates results for iterative methods for solving nonlinear ill-posed problems. In *Surveys on solution methods for inverse problems*, 7-34, Springer-Verlag, Vienna, 2000.



K. Eriksson, D. Estep, P. Hansbo and C. Johnson. *Computational differential equations*. Cambridge University Press, Cambridge, 1996.



K. Eriksson and C. Johnson. An adaptive finite element method for linear elliptic problems, *Math. Comput.*, 50(1988), 361-383.



K. Eriksson and C. Johnson. Adaptive finite element methods for parabolic problems. I. A linear model problem, *SIAM J. Numer. Anal.*, 28 (1991), 1, 43-77.



R. E. Ewing. A posteriori error estimation, *Comput. Meth. Appl. Mech. Engrg.*, 82(1990),1-3, 59-72.



R. E. Ewing, R. D. Lazarov and J. Wang, Superconvergence of the velocity along the Gauss lines in mixed finite element methods, *SIAM J. Numer. Anal.*, 18(1991), 1015-1029.



R. S. Falk. Error estimates for the approximation of a class of variational inequalities, *Math. Comput.*, 28(1974), 963-971.



M. Feistauer. *Mathematical methods in fluid dynamics*. Longman, Harlow, 1993.



W. Fenchel. On the conjugate convex functions, *Canad. J. Math.* 1(1949), 73-77.



Fraeys B. de Veubeke. Displacement and equilibrium models in the finite element methods, in *Stress analysis* (O. C. Zienkiewicz and G. S. Holister, eds.), Wiley and Sons, New-York, 1965, 145-197.



Fraeys B. de Veubeke. A conforming finite element for plate bending, *Internat. J. Solids and Structures* 4(1968), 95-108.



G. E. Forsythe, M. A. Malcolm and C. B. Moler. *Computer methods for mathematical computations*, Prentice-Hall, Englewood Cliffs, New York, 1977.



A. Friedman. *Variational principles and free-boundary problems*. Wiley and Sons, New York, 1982.



D. A. Field. Laplacian smoothing and Delaney triangulations, *Commun. Appl. Numer. Methods*, 4(1988), 709-712.



L. P. Franca, J. Karam Filho, A. F. D. Loula, and R. Stenberg, *A convergence analysis of a stabilized method for the Stokes flow*, *Mat. Apl. Comput.* 10 (1991),1, pp. 19–26.



M. Frolov, P. Neittaanmäki and S. Repin. On the reliability, effectivity and robustness of a posteriori error estimation methods. Reports of the Department of Mathematical Information Technology of the University of Jyväskylä, No. B14/2002.



M. Frolov, P. Neittaanmäki and S. Repin. On computational properties of a posteriori error estimates based upon the method of duality error majorants (to appear in Proc. 5th European Conference on Numerical Mathematics and applications, Praha, 2003).



M. Fuchs and G.A. Seregin. *Variational methods for problems from plasticity theory and for generalized Newtonian fluids*. Lect. Notes in Mathematics 1749, Springer-Verlag, Berlin (2000).



H. Gaevskii, H., K. Gröger, and K. Zacharias. *Nichtlineare Operator-gleichungen and Operatordifferentialgleichungen*. Akademie-Verlag, Berlin, 1974.



B. G. Galerkin. Beams and plates. Series in some questions of elastic equilibrium of beams and plates. *Vestnik Ingenerov*, 19(1915), 897-908 (in Russian).



L. Gallimard, P. Ladevèze and J. P. Pelle. Error estimation and adaptivity in elastoplasticity, *Int. J. Numer. Meth. Engrg.*, 39(1996), 189-217.



E.G. Geisler, A.A. Tal and D.P. Garg. On a-posteriori error bounds for the solution of ordinary nonlinear differential equations. In *Computers and mathematics with applications*, 407-416, Pergamon, Oxford, 1976.



C. I. Goldstain. Variational crimes and L^∞ error estimates in the finite element method, *Math. Comput.*, 35(1980), 152, 1131-1157.



D. Gilbarg and N.S. Trudinger. *Elliptic partial differential equations of second order*. Springer-Verlag, Berlin, 1977.



V. Girault and P. A. Raviart. "*Finite element approximation of the Navier–Stokes equations*". Springer-Verlag, Berlin, 1986.



R. Glowinski. *Numerical methods for nonlinear variational problems*. Springer-Verlag, New-York, 1982.



R. Glowinski and O. Pironneau. Numerical Methods for the First Biharmonic Equation and for the Two-Dimensional Stokes Problem, *SIAM Review*, (21)1999, 2, 167-212.



R. Glovinski, J.-L.Lions, R. Trémolierès. *Analyse numérique des inéquations variationnelles*. Dunod, Paris, 1976.



R. Glowinski, T.-W. Pan and J. Periaux. A fictitious domain method for Dirichlet problem and applications, *Comput. Methods Appl. Mech. Engrg.*, 111(1994), 283-303.



W. Han. Finite element analysis of a holonomic elasto-plastic problem, *Numer. Math.*, 60(1992), 493-508.



W. Han. A posteriori error analysis for linearization of nonlinear elliptic problems and their discretization, *Math. Meth. Appl. Sci.* , 17(1994), 487-508.



W. Han and D. C. Reddy. On the finite element method for mixed variational inequalities arising in elastoplasticity, *SIAM J. Numer. Anal.*, 32(1995), 6, 1776-1807.



W. Han and D. Reddy. Qualitative and numerical analysis of quasi-static problems in elastoplasticity, *SIAM J. Numer. Anal.*, 34(1997), 1, 143-177.



W. Han, S. Jensen and I. Shimansky. The Kačanov method for some nonlinear problems, *Applied Numerical Mathematics*, 24(1997), 57-79.



P. Hansbo. Generalized Laplacian smoothing of unstructured grids, *Commun. Numer. Methods Eng.*, 11(1995), 455-464.



Y. Hayashi. On a posteriori error estimation in the numerical solution of systems of ordinary differential equations. *Hiroshima Math. J* 9(1979), no, 1, 201-243.



R.Hill. *The mathematical theory of plasticity*. Oxford, 1983



I. Hlaváček and M. Křížek. On a superconvergence finite element scheme for elliptic systems. I. Dirichlet boundary conditions. *Aplikace Matematiky*, 32(1987), No.2, 131-154.



R. H. W. Hoppe. Mortar finite elements in R^3 , *East-West J. Numer. Math.*, 7(1999), 3, 159-173.



R. H. W. Hoppe and R. Kornhuber. Adaptive multilevel methods for obstacle problems. *SIAM J. Numer. Anal.*, 31(1994), 301-323.



P. Houston, R. Rannacher, E. Süli. A posteriori error analysis for stabilised finite element approximations of transport problems, *Comput. Methods Appl. Mech. Engrg.* 190, 1483-1508, 2000.



A. A. Iljushin. *Plasticite. (Deformations elastico-plastiques)*. (Eyrolles, Paris, 1956).



A.D. Ioffe and V.M. Tikhomirov. *Theory of extremal problems. Studies in Mathematics and its Applications, 6*. North-Holland, Amsterdam-New York, 1979



J.L. Jensen. Sur les fonctions convexes et les inegalités entre les valeurs moyennes, *Acta Math.*, 30(1906), 175-193.



Jin Qi-nian. Error estimates of some Newton-type methods for solving nonlinear inverse problems in Hilbert scales, *Inverse Problems*, 16(2000), 187-197



C. Johnson. *Numerical Solution of Partial Differential Equations by the Finite Element Methods*, Cambridge University Press, Cambridge, 1987.



C. Johnson. On Finite Element Methods for Plasticity Problems.-
Numer. Math., 26(1976), 79-84.



C. Johnson. Adaptive finite element methods for diffusion and convection problems, *Comput. Methods Appl. Mech. Engrg.*, 82(1990), 301-322.



C. Johnson and P. Hansbo. Adaptive finite elements in computational mechanics, *Comput. Methods Appl. Mech. Engrg.* 101(1992), 143-181.



C. Johnson and B. Mercier. Some equilibrium finite element methods two-dimensional elasticity problems, *Numer. Math.*, 30(1978), 101-116.



C. Johnson and R. Rannacher. On error control in computational fluid dynamics (CFD) Preprint 1994-07, Chalmers University of Technology, Göteborg.



C. Johnson, R. Rannacher and M. Boman. Numerics in hydrodynamic stability. Towards error control in CFD, *SIAM J. Numer. Anal.*, 32(1995), 1058-1079.



C. Johnson, R. Rannacher and M. Boman. *On transition to turbulence and error control in CFD*, Preprint 1994–26, Chalmers University of Technology, Göteborg.



C. Johnson and A. Szepessy. Adaptive finite element methods for conservation laws based on a posteriori error estimates, *Commun. Pure and Appl. Math.*, vol. XLVIII (1995), 199-234.



B.-O. Heimsund, X.-C. Tai and J. Wang. Superconvergence for the gradient of finite element approximations by L^2 projections, *SIAM J. Numer. Anal.*, 40(2002), 4, 1263-1280.



H. Kavarada. *Numerical problems for free surface problems by means of penalty*. Lecture Notes in Mathematics, No 704, Springer-Verlag, 1979.



H. Kardestuncer and D. H. Norrie ed. *Finite Element Handbook*. McGraw–Hill, New York, 1987.



D. W. Kelly. The self equilibration of residuals and complementary error estimates in the finite element method, *Internat. J. Numer. Meth. Engrg.* 20(1984) 1491-1506.



D. Kinderlehrer and G. Stampacchia. *An introduction to variational inequalities and their applications*. Academic Press, New York, 1980.



D. W. Kelly, J. R. Gago, O. C. Zienkiewicz and I. Babuška. A posteriori error analysis and adaptive processes in the finite element method. Part I - error analysis, *Internat. J. Numer. Meth. Engrg.*, 19(1983), 1593-1619.



R. V. Kohn. The relaxation of a double-well energy, *Continuum Mech. Thermodynamics*, 3(1991), 3, 193-236.



A. N. Kolmogorov and S. V. Fomin. *Introductory real analysis*. Dover Publications, Inc., New York, 1975.



R. Kornhuber. A posteriori error estimates for elliptic variational inequalities, *Comput. Math. Appl.*, 31(1996), 49-60.



S. Korotov, P. Neittaanmaki and S. Repin. A posteriori error estimation of goal-oriented quantities by the superconvergence patch recovery, *J. Numer. Math.* 11 (2003), 1, 33-59.



M. Křížek, P. Neittaanmäki and R Stenberg eds. *Finite Element Methods. Superconvergence, Post-Processing and A Posteriori Error Estimates*. Lecture Notes in Pure and Applied Mathematics, Vol. 196, Marcel Dekker, New York, 1998.



M. Křížek and P. Neittaanmäki. Superconvergence phenomenon in the finite element method arising from averaging of gradients *Numer. Math.*, 45(1984), 105-116.



P. Ladeveze and D. Leguillon. Error estimate procedure in the finite element method and applications, *SIAM J. Numer. Anal.*, 20(1983), 485-509.



P. Ladeveze, J.-P. Pelle and Ph. Rougeot. Error estimation and mesh optimization for classical finite elements, *Engineering Computations*, 8(1991), 69-80.



P. Ladev eze and Ph. Rougeot. New advances on a posteriori error on constitutive relation in f.e. analysis, *Comput. Methods Appl. Mech. Engrg.*, 150(1997), 239-249.



O. A. Ladyzhenskaya. *Mathematical problems in the dynamics of a viscous incompressible fluid*. Nauka, Moscow, 1970 (in Russian).



O. A. Ladyzhenskaya, *The boundary value problems of mathematical physics*. Springer-Verlag, New York, 1985.



O. A. Ladyzhenskaya, On modified Navier–Stokes equations for large velocity gradients, *Zapiski Nauchnykh Seminarov LOMI*, 7(1968), 126-154.



O. A. Lady zenskaja and V. A.; Solonnikov. Some problems of vector analysis, and generalized formulations of boundary value problems for the Navier-Stokes equation, *Zap. Nauchn. Sem. Leningrad. Otdel. Mat. Inst. Steklov. (LOMI)*, 59(1976), 81–116, 256 (Russian).



O. A. Ladyzhenskaya and N. N. Uraltseva. *Linear and Quasilinear Elliptic equations*, Academic Press, New York, 1968.



R. D. Lazarov. Superconvergence of the gradient for triangular and tetrahedral finite elements of a solution of linear problems in elasticity theory, *Computational Processes and Systems*, 6(1988), 180-191 (in Russian).



A.S. Leonov. Some a posteriori stopping rules for iterative methods for solving linear ill-posed problems, *Comput. Math. Math. Phys.*, 34(1994), 1, 121-126.



J.-L. Lions and E. Magenes. Problèmes aux limites non homogènes et applications. Dunod, Paris, 1968.



R. Löhner, K. Morgan and O. C. Zienkiewicz. Adaptive grid refinement for the compressible Euler equations. In *Accuracy Estimates and Adaptive Refinements in Finite Element Computations, I*. Babuška, O. C. Zienkiewicz, J. Gago and E. R. de Oliveira eds., Wiley and Sons, 1986, 281-297.



J. Malek, J. Nečas, J. Rokuta and M. Ružička. *Weak and measure valued solutions to evolution partial differential equations. Applied*

Mathematic and Mathematical Computation vol 13., Chapman and Hall, 1996.



J. Medina, M. Picasso and J. Rappaz. Error estimates and adaptive finite elements for nonlinear diffusion-convection problems. Ecole Polytechnique Federale de Lausanne, Preprint CH-1015, 1995.



P. Meyer. A unifying theorem on Newton's method, *Numer. Funct. Anal. Optim.* 13(1992), no. 5-6, 463-473.



S. G. Mikhlin. *Variational methods in mathematical physics.* Pergamon, Oxford, 1964.



S. G. Mikhlin. *Error Analysis in Numerical Processes* Wiley and Sons, Chicester–New York, 1991.



S. G. Mikhlin. *Constants in some inequalities of analysis.* Wiley and Sons, Chicester–New York, 1986.



P. Mosolov and V. Myasnikov. *Mechanics of rigid plastic bodies*, Nauka, Moscow, 1981 (in Russian).



A. Muzalevsky and S. Repin. On two-sided error estimates for approximate solutions of problems in the linear theory of elasticity, *Russian J. Numer. Anal. Math. Modelling*, 18(2003), 1, 65-85.



P. Neittaanmäki and M. Křížek. On $O(h^4)$ superconvergence of piecewise bilinear FE-approximations. In Teubner-Texte Math. 107, Teubner, 1988, 250-255.



P. Neittaanmäki and S. Repin. A posteriori error estimates for boundary-value problems related to the biharmonic operator, *East-West J. Numer. Math.*, 9(2001), 2, 157-178.



P. Neittaanmäki and S. Repin, *Reliable methods for computer simulation, Error control and a posteriori estimates*, Elsevier, New York, 2004.



J. T. Oden, L. Demkowicz, T. Strouboulis and P. Devloo. Adaptive methods for problems in solid and fluid mechanics. In *Accuracy Estimates and Adaptive Refinements in Finite Element Computations, I*. Babuška, O. C. Zienkiewicz, J. Gago and E. R. de Oliveira eds., Wiley and Sons, 1996.



J. T. Oden, L. Demkowicz, W. Rachowicz, and T. A. Westermann. Towards a universal $h - p$ adaptive finite element strategy. Part 2. A posteriori error estimation, *Comput. Methods Appl. Mech. Engrg.*, 77 (1989) 113-180.



J. T. Oden, S. Prudhomme. Goal-oriented error estimation and adaptivity for the finite element method, *Comput. Math. Appl.*, 41, 735-756, 2001.



J. T. Oden and J. N. Reddy. *Mathematical theory of finite elements*. Wiley and Sons, 1976.



J. T. Oden, W. Wu and M. Ainsworth. An a posteriori error estimate for finite element approximations of the Navier–Stokes equations, *Comput. Methods. Appl. Mech. Engrg.*, 111(1994), 185-202.



L. A. Oganessian and L. A. Ruchovet. An investigation of the rate of convergence of variation-difference schemes for second order elliptic equations in a two-dimensional region with smooth boundary, *Z. Vychisl. Mat. i Mat. Fiz.*, 9(1969), 1102–1120, (Russian).



A. Ostrowski. Les estimations des erreurs a posteriori dans les procédés itératifs, *C.R. Acad.Sci. Paris Sér. A-B*, 275(1972), A275-A278.



D.R.J.Owen and E.Hinton. *Finite Elements in Plasticity Theory: Theory and Practice*. Prinerdge Press, Swansea, 1980.



C. Padra, A posteriori error estimators for nonconforming approximation of some quasi-newtonian flows, *SIAM J. Numer. Anal.*, 34(1997), 1600-1615.



J. Peraire and A. T. Patera. Bounds for linear-functional outputs of coercive partial differential equations: Local indicators and adaptive refinement. In *Advances in Adaptive Computational Methods in Mechanics*, Ed. P. Ladevéze and J. T. Oden, Elsevier, New York, 1998, 199-228.



F. Potra. Sharp error bounds for a class of Newton-like methods, *Libertas Math.* 5(1985), 71-84



J. Pousin and J. Rappaz. Consistance, stabilité, erreurs a priori et a posteriori pour des problemes non linéaires. *C. R. Acad. Sci. Paris* 312(1991), 699-703.



J. Pousin and J. Rappaz. Consistency, stability, a priori and a posteriori errors for Petrov-Galerkin methods applied to nonlinear problems. *Numer. Math.* 69(1994), 213-231.



W. Prager and J. L. Synge. Approximation in elasticity based on the concept of function space, *Quart. Appl. Math.* 5(1947), 241-269.



R. Rannacher. Zur L^∞ -Konvergenz linear finite elemente beim Dirichlet Problem, *Math. Z.*, 149(1976), 69-77.



R. Rannacher. On nonconforming and mixed finite element method for plate bending problems, the linear case, *R.A.I.R.O. Anal. Numer.*, 13(1979), 4, 369-387.



R. Rannacher and R. Scott. Some optimal error estimates for piecewise linear finite element approximations, *Math. Comput.*, 38(1982), 158, 437-445.



R. Rannacher and F. T. Suttmeier. A feed-back approach to error control in finite element methods: application to linear elasticity. IWR, Preprint 96-42(SFB 359), Heidelberg 1996.



R. Rannacher and F.T. Suttmeier. A posteriori error control and mesh adaptation for F.E. models in elasticity and elasto-plasticity. In *Advances in Adaptive Computational Methods in Mechanics*, Eds. P. Ladevéze and J. T. Oden, Elsevier, New York, 1998, 275-292.



W. C. Rheinboldt. On a theory of mesh-refinement processes, *SIAM J. Numer. Anal.*, 17(1980), 766-778.



S. Repin. A posteriori estimates for approximate solutions of variational problems with strongly convex functionals, *Problems of Mathematical Analysis*, 17 (1997), 199-226. (in Russian, translated in *Journal of Mathematical Sciences*, 97(1999), 4, 4311-4328).



S. Repin. A posteriori estimates of the accuracy of variational methods for problems with nonconvex functionals, *Algebra i Analiz*, 11(1999), 4, 151-182 (in Russian, translated in *St.-Petersburg Mathematical Journal*, 11(2000), 4, 651-672).



S. Repin. A posteriori error estimation for nonlinear variational problems by duality theory. *Zapiski Nauchnykh Seminarov POMI*, 243(1997) 201-214.



S. Repin. A posteriori error estimation for variational problems with power growth functionals based on duality theory, *Zapiski Nauchnykh Seminarov POMI*, 249(1997), 244-255.



S. Repin. A posteriori error estimates for approximate solutions of variational problems. In *Proceedings of 2nd European Conference on Numerical Mathematics and Advanced Applications, Heidelberg, 1997*, 524-531, World Sci. Publishing, River Edge, New York, 1998.



S. Repin. A unified approach to a posteriori error estimation based on duality error majorants, *Mathematics and Computers in Simulation*, 50(1999), 313-329.



S. Repin. A posteriori error estimation for variational problems with uniformly convex functionals, *Math. Comput.*, 69(230), 2000, 481-500.



S. Repin. Two-sided estimates for deviation from an exact solution to uniformly elliptic equation. *Trudi St.-Petersburg Math. Society*, 9(2001), 148-179 (in Russian, translated in *American Mathematical Translations Series 2*, 9(2003))



S.Repin. Estimates of deviation from exact solutions of initial-boundary value problems for the heat equation, *Rend. Mat. Acc. Lincei*, 13(2002), 121-133.



S. Repin. Estimates of deviations from exact solutions of elliptic variational inequalities, *Zapiski Nauchn. Semin. V.A. Steklov Mathematical Institute in St.-Petersburg (POMI)*, 271(2000), 188-203.



S. Repin. Aposteriori estimates for the Stokes problem, *Journal of Math. Sciences* 109 (2002), 5, 1950-1964.



S. Repin. Estimates of deviations for generalized Newtonian fluids, *Zapiski Nauchn. Semin. V.A. Steklov Mathematical Institute in St.-Petersburg (POMI)*, 288(2002), 178-203.



S. Repin. Estimates for errors in two-dimensional models of elasticity theory, *J. Math. Sci. (New York)*, 106 (2001), no. 3, 3027-3041.



S. Repin, *Functional Approach to Locally Based A Posteriori Error Estimates for Elliptic and Parabolic Problems*, Proc. 6th European

Conference on Numerical Mathematics and Advanced Applications (Santiago de Compostela) 2005, Springer, Berlin, 2006, pp. 133-148.



S. Repin, *Local a posteriori estimates for the Stokes problem*, Zap. Nauchn. Sem. S.-Peterburg, Otdel. Mat. Inst. Steklov. (POMI) 318 (2004), pp. 233–245.



S. Repin and M. Frolov. A posteriori error estimates for approximate solutions of elliptic boundary value problems, *Zh. Vychisl. Mat. Mat. Fiz.*, 42(2002), 12, 1774–1787 (Russian).



S. Repin, S. Sauter and A. Smolianski. A posteriori error estimation for the Dirichlet problem with account of the error in the approximation of boundary conditions, *Computing*, 70(2003), 205-233.



S. Repin, S. Sauter and A. Smolianski. Duality Based A Posteriori Error estimator for the Dirichlet Problem, *Proc.Appl. Math. Mech.*, 2 (2003), 513-514.



S. Repin, S. Sauter and A. Smolianski. A posteriori error estimation of dimension reduction errors (to appear in Proc. 5th European

Conference on Numerical Mathematics and applications, Praha,2003).



S. Repin, S. Sauter and A. Smolianski, *A posteriori estimation of dimension reduction errors for elliptic problems in thin domains*, *SIAM J. Numer. Anal.*, 42 (2004), 4, pp. 1435–1451.



S. Repin, S. Sauter and A. Smolianski, *A Posteriori Control of Dimension Reduction Errors on Long Domains*, *Proceedings in Applied Mathematics and Mechanics*, 4, 1, 714–715 (2004).



S. I. Repin and L. S. Xanthis. A posteriori error estimation for elasto-plastic problems based on duality theory, *Comput. Methods Appl. Mech. Engrg.*, 138(1996), 317-339.



S. I. Repin and L. S. Xanthis. A posteriori error estimation for nonlinear variational problems, *Comptes Rendus de l'Académie des Sciences, Mathématique*, 324(1997), 1169-1174.



W. Ritz. Über eine neue Methode zur Lözing gewisser Variationsprobleme der Mathematischen Physics, *J. Reine Angew. Math.*, 135(1909), 1-61.



R. T. Rockafellar. *Convex analysis*, Princeton Univ. Press, 1970.



R. Rodrigues. Some remarks on Zienkiewicz–Zhu estimator, *Numer. Methods for PDE*, 10(1994), 625-635.



A. H. Schatz. Pointwise error estimates and asymptotic error expansion inequalities for the finite element method on irregular grids: Part 1: global estimates, *Math. Comput.*, 67(1998), 223, 877-899.



A. H. Schatz and L. B. Wahlbin. Interior maximum norm estimates for finite element methods, *Math. Comput.*, 31(1977), 414-442.



A. H. Schatz and L. B. Wahlbin. Maximum norm estimates for finite element method on plane polygonal domains. Part 1, *Math. Comput.*, 32(1978), 73-109.



A. H. Schatz and L. B. Wahlbin. Maximum norm estimates for finite element method on plane polygonal domains. Part 2, refinements, *Math. Comput.*, 33(1979), 465-492.



C. Schwab. A-posteriori modelling error estimation for hierarchic plate models, *Numer. Math.*, 74(1996), 221-259.



A. H. Schatz and J. Wang. Some new error estimates for Ritz–Galerkin methods with minimal regularity assumptions, *Math. Comput.*, 65(1996), 213, 19-27.



M. Schultz and O. Steinbach. A new a posteriori error estimator in adaptive direct boundary element methods: the Dirichlet problem, *Calcolo*, 37(2000), 79-96.



M. Schulz and W. L. Wendland. A general approach to a posteriori error estimates for strictly monotone and Lipschitz continuous nonlinear operators illustrated in elasto-plasticity. In *Proceedings of 2nd European Conference on Numerical Mathematics and Advanced Applications, Heidelberg, 1997*, World Sci. Publishing, River Edge, New York, 1998, 572-579.



M. Schultz and W. Wendland. *On an adaptive finite element method for elasto-plastic deformation with hardening*. Preprint 1997/39, SFB 259, Stuttgart, 1997.



S. L. Sobolev. *Some Applications of Functional Analysis in Mathematical Physics*, Izdat. Leningrad. Gos. Univ., Leningrad, 1955 (in Russian)

translated in *Translation of Mathematical Monographs, Volume 90* American Mathematical Society, Providence, RI, 1991).



E. Stein and S. Ohnibus. Coupled model- and solution-adaptivity in the finite element method, *Comput. Methods Appl. Mech. Engrg.*, 150(1997), 327-350.



E. Stein, F.J. Bartold, S. Ohnibus and M. Schmidt. Adaptive finite elements in elastoplasticity with mechanical error indicators and Neumann-type estimators. In *Advances in Adaptive Computational Methods in Mechanics*, Ed. P. Ladevéze and J. T. Oden, Elsevier, New York, 1998, 81-100.



G. Strang and G. Fix. *An analysis of the finite element method*. Prentice Hall, Englewood Cliffs, 1973.



T. Strouboulis and J. T. Oden. A posteriori error estimation of finite element approximations in fluid mechanics, *Comput. Meth. Appl. Mech. Engrg.*, 78(1990), 201-242.

-  J. L. Synge. The hypercircle method. In *Studies in numerical analysis (papers in honour of Cornelius Lanczos on the occasion of his 80th birthday)*, 201-217. Academic Press, London, 1974.
-  R. Temam, *Navier–Stokes equations. Theory and numerical analysis*, Studies in Mathematics and its Applications, 2, North-Holland, Amsterdam, 1979.
-  K. Tsuruta and K. Ohmori. A posteriori error estimation for Volterra integro-differential equations, *Mem. Numer. Math. No. 3* 1976, 33-47.
-  N.N. Uraltseva. On the regularity of solutions to variational inequalities *Uspekhi Mat. Nauk*, 42(1987), 6, 151-174 (in Russian).
-  R.S. Varga. *Matrix Iterative Analysis*. Prentice–Hall, Englewood Cliffs, New Jersey, 1962.
-  A. Ve eser. Efficient and reliable a posteriori error estimators for elliptic obstacle problems, *SIAM J. Numer. Anal.*, 39(2001), 146-167.
-  R. Verfürth. A posteriori error estimators for the Stokes equations, *Numer. Math.*, 55(1989), 309-326.



R. Verfürth. A posteriori error estimates for nonlinear problems. Finite element discretisations of elliptic equations, *Math. Comput.* 62(1994) 445-475.



R. Verfürth. *A review of a posteriori error estimation and adaptive mesh-refinement techniques* Wiley and Sons, Teubner, New-York, 1996.



R. Verfürth. A posteriori error estimates for nonlinear problems. $L^r(0, T; L^p(\Omega))$ -error estimates for finite element discretizations of parabolic equations, *Math. Comput.*, 67(1998), 224, 1335-1360.



A. Vesser. Efficient and reliable a posteriori error estimators for elliptic obstacle problems, *SIAM J. Numer. Anal.*, 39(2001), 1, 146-167.



M. I. Vishik . The method of orthogonal projections for selfadjoint equations , *Sov. Math. Dokl* , 1947, 56.



L. B. Wahlbin. *Superconvergence in Galerkin Finite Element Methods*, Lecture Notes in Mathematics, No 1605, Springer-Verlag, 1995.



J. Wang. Superconvergence analysis of finite element solutions by the least-squares surface fitting on irregular meshes for smooth problems, *J. Math. Study*, 33(2000), 3, 229-243.



H. Weil. The method of orthogonal projections in potential theory, *Duke Math. J.*, 7(1940), 411-444.



J. R. Whiteman and G. Goodsell. A survey of gradient superconvergence for finite approximations to second order elliptic problems on triangular and tetrahedral meshes. In *The Mathematics of Finite Elements and Applications VII*, Ed. J. R. Whiteman, Academic Press, 1991, 55-74.



W. Wunderlich, H. Gramer and G. Steinl. An adaptive finite element approach in associated and non-associated plasticity considering localization phenomena. In *Advances in Adaptive Computational Methods in Mechanics*, Ed. P. Ladevéze and J. T. Oden, 293-308, Elsevier, New York, 1998.



On classes of summable functions and their Foutier series, *Proc. Roy. Soc. Ser. A*, 87(1912), 225-229.



S. Zaremba. Sur un probleme toujours possible comprenant, a titre de cas particuliers, le probleme de Dirichlet et celui de Neumann, *J. math. pures et appl.* 6(1927),2, 127-163.



E. Zeidler. *Nonlinear functional analysis and its applications. I. Fixed-point theorems.* Springer-Verlag, New York, 1986,



O. C. Zienkiewicz. Achievements and some unsolved problems of the finite element method, *Int. J. Numer. Meth. Engrg.*, 47(2000), 9-28.



O. C. Zienkiewicz, B. Boroomand and J.Z. Zhu. Recovery procedures in error estimation and adaptivity: Adaptivity in linear problems. In *Advances in Adaptive Computational Methods in Mechanics*, Ed. P. Ladevéze and J. T. Oden, Elsevier, New York, 1998, 3-23.



O. C. Zienkiewicz and K. Morgan. *Finite elements and approximation.* Wiley and Sons, New York, 1983.



O.C.Zienkiewicz, S.Valliappan, I.P.King. Elasto-plastic solutions of engineering problems, initial stress finite element approach.- *Int. J. Numerical Methods Engrg.*, 1969, 1, 75-100.



O. C. Zienkiewicz and J. Z. Zhu. A simple error estimator and adaptive procedure for practical engineering analysis, *Internat. J. Numer. Meth. Engrg.*, 24(1987) 337-357.



O. C. Zienkiewicz and J. Z. Zhu. Adaptive techniques in the finite element method, *Commun. Appl. Numer. Methods*, 4(1988), 197-204.



O. C. Zienkiewicz and J. Z. Zhu. The superconvergent patch recovery and a posteriori error estimates. Part 1: The recovery technique, *Int. J. Numer. Meth. Engrg.*, 33(1992), 1331-1364.



J. Z. Zhu. A posteriori error estimation - the relationship between different procedures, *Comput. Methods Appl. Mech. Engrg.*, 150(1997), 411-422.



M. Zlámal. Some superconvergence results in the finite element method. Mathematical aspects of finite element methods. In *Proc. Conf., Consiglio Naz. delle Ricerche (C.N.R.), Rome, 1975. Lecture Notes in Math., Vol. 606, Springer-Verlag, Berlin, 1977, 353-362.*