

Charles University in Prague
Faculty of Mathematics and Physics

MASTER THESIS



Lukáš Korous

Adaptive hp Discontinuous Galerkin Method for Nonstationary Compressible Euler Equations

Department of Numerical Mathematics

Supervisor: prof. RNDr. Miloslav Feistauer, DrSc., dr. h. c.

Study branch: Numerical and Computational Mathematics

2011

I would like to thank prof. RNDr. Miloslav Feistauer, DrSc., dr. h. c. for his guidance, attentiveness and valuable advice I received from him. Also I would like to thank doc. RNDr. Pavel Šolín, Ph.D for the priceless advices, and enjoyable collaboration on the Finite element software Hermes.

I confirm having prepared the master thesis by my own, and having listed all used sources of information in the bibliography. I agree with lending the master thesis.

In Prague

Lukáš Korous

Contents

Introduction	5
1 Compressible flow and the Euler equations	6
1.1 Equations describing the flow	6
1.1.1 Description of the flow	6
1.1.2 The Transport theorem	8
1.1.3 The continuity equation	9
1.1.4 The equations of motion	11
1.1.5 The Navier-Stokes equations	12
1.1.6 The energy equation	13
1.1.7 Thermodynamical relations	14
1.1.8 Entropy and the second law of thermodynamics	15
1.2 Euler equations and their properties	18
1.2.1 Dimensionless Euler equations	18
1.2.2 Conservative form of the Euler equations	19
1.2.3 Rotational invariance	22
1.2.4 Diagonalization of the Jacobi matrix	22
2 Computational methods in fluid dynamics	24
2.1 Triangulation	24
2.2 Finite Volume method	26
2.2.1 Scheme derivation	26
2.2.2 Properties of the numerical flux	26
2.2.3 Example	27
2.3 Finite Element Method	28
2.3.1 Example	29
2.4 Discontinuous Galerkin method	29
2.4.1 Derivation of the DGFEM	30
2.4.2 Example	32
3 Adaptive algorithms	33
3.1 p -adaptivity	34
3.1.1 Higher-order finite element space	34

3.1.2	Hierarchic shape functions	34
3.1.3	Higher-order numerical quadrature	38
3.2	Hanging nodes	40
3.2.1	Hanging nodes and irregularity rules	40
3.3	Data structures for <i>hp</i> -adaptivity	44
3.3.1	Node and element structures	44
3.3.2	Eliminating neighbor search by hashing	45
3.3.3	Finding all neighbors of an element	47
3.3.4	Determining node type	47
3.4	Automatic <i>hp</i> -adaptivity	49
3.4.1	Reference solution	49
3.4.2	Outline of the algorithm	49
3.4.3	Selecting optimal <i>hp</i> -refinement	51
3.4.4	Automatic <i>hp</i> -adaptivity for time-dependent problems	56
4	Dicontinuous Galerkin discretization of the Euler equations	60
4.1	Numerical Fluxes	60
4.2	Time discretization	62
4.2.1	Linearization	63
4.3	Boundary conditions	64
4.3.1	Semi-implicit method for the time integration	65
4.3.2	Stability and the CFL condition	65
4.4	Shock capturing	65

Název práce: Adaptivní hp nespojitá Galerkinova metoda pro nestacionární stlačitelné Eulerovy rovnice

Autor: Lukáš Korous

Katedra (ústav): Katedra numerické matematiky

Vedoucí bakalářské práce: prof. RNDr. Miloslav Feistauer, DrSc., dr. h. c.

e-mail vedoucího: feist@karlin.mff.cuni.cz

Abstrakt: Stlačitelné Eulerovy rovnice popisují pohyb stlačitelných nevazkých tekutin. Používají se v mnoha oblastech leteckého, automobilového a jaderného inženýrství, chemie, ekologie, klimatologie, i jinde. Matematicky, stlačitelné Eulerovy rovnice představují hyperbolický systém skládající se z několika nelineárních parciálních diferenciálních rovnic (zákony zachování). Tyto rovnice jsou řešeny nejčastěji pomocí metody konečných objemů (MKO), a metody konečných prvků (MKP) nízkého řádu. Nicméně, oba tyto přístupy nedosahují vyššího řádu přesnosti, a navíc je dobře známo, že konformní metoda konečných prvků není optimální nástroj pro diskretizaci rovnic prvního řádu. Nejnadějnější přístup k přibližnému řešení stlačitelných Eulerových rovnic je nespojitá Galerkinova metoda, která kombinuje stabilitu MKO s vynikajícími aproximačními vlastnostmi MKP vyššího řádu. Cílem této diplomové práce byl vývoj, implementace a testování nových algoritmů pro adaptivní řešení nestacionárních stlačitelných Eulerových rovnic na základě vyššího řádu nespojité Galerkinovy metody (hp-DG). Základem pro nové metody byly nespojitá Galerkinova metoda a časoprostorové hp-MKP algoritmy na dynamických sítích pro nestacionární problémy druhého řádu. Nové algoritmy byly implementovány a testovány v rámci open source knihovny Hermes, a vztahují se na vybrané problémy transonického proudění, a proudění s nízkým Machovým číslem v klimatologii.

Klíčová slova: numerické simulace, metoda konečných prvků, Eulerovy rovnice, hp-adaptivita, nespojitá Galerkinova metoda

Title: Adaptive hp Discontinuous Galerkin Method for Nonstationary Compressible Euler Equations

Author: Lukáš Korous

Department: Department of Numerical Mathematics

Supervisor: prof. RNDr. Miloslav Feistauer, DrSc., dr. h. c.

Supervisor's e-mail address: feist@karlin.mff.cuni.cz

Abstract: The compressible Euler equations describe the motion of compressible inviscid fluids. They are used in many areas ranging from aerospace, automotive, and nuclear engineering to chemistry, ecology, climatology, and others. Mathematically, the compressible Euler equations represent a hyperbolic system consisting

of several nonlinear partial differential equations (conservation laws). These equations are solved most frequently by means of Finite Volume Methods (FVM) and low-order Finite Element Methods (FEM). However, both these approaches are lacking higher order accuracy and moreover, it is well known that conforming FEM is not the optimal tool for the discretization of first-order equations. The most promising approach to the approximate solution of the compressible Euler equations is the discontinuous Galerkin method that combines the stability of FVM, with excellent approximation properties of higher-order FEM. The objective of this Master Thesis was to develop, implement and test new adaptive algorithms for the nonstationary compressible Euler equations based on higher-order discontinuous Galerkin (hp-DG) methods. The basis for the new methods will be the discontinuous Galerkin methods and space-time adaptive hp-FEM algorithms on dynamical meshes for nonstationary second-order problems. The new algorithms will be implemented and tested in the framework of the open source library Hermes, and they will be applied to selected problems of transonic flow and low Mach number flows in climatology.

Keywords: numerical simulation, finite element method, Euler equations, hp-adaptivity, discontinuous Galerkin method

Introduction

Discontinuous Galerkin methods (DG methods) in mathematics form a class of numerical methods for solving partial differential equations. They combine features of the finite element and the finite volume framework and have been successfully applied to hyperbolic, elliptic and parabolic problems arising from a wide range of applications. DG methods have in particular received considerable interest for problems with a dominant first-order part, e.g. in electrodynamics, fluid mechanics and plasma physics.

Discontinuous Galerkin methods were first proposed and analyzed in the early 1970s as a technique to numerically solve partial differential equations. In 1973 Reed and Hill introduced a DG method to solve the hyperbolic neutron transport equation.

The origin of the DG method for elliptic problems cannot be traced back to a single publication as features such as jump penalization in the modern sense were developed gradually. However, among the early influential contributors were Babuška, J.-L. Lions, Nitsche and Zlamal. Interestingly DG methods for elliptic problems were already developed in a paper by Baker in the setting of 4th order equations in 1977. A more complete account of the historical development and an introduction to DG methods for elliptic problems is given in a publication by Arnold, Brezzi, Cockburn and Marini. A number of research directions and challenges on DG methods are collected in the proceedings volume edited by Cockburn, Karniadakis and Shu.

The discontinuous Galerkin (DG) methods are locally conservative, stable, and high-order accurate methods which can easily handle complex geometries, irregular meshes with hanging nodes, and approximations that have polynomials of different degrees in different elements. These properties, which render them ideal to be used with hp-adaptive strategies, not only have brought these methods into the main stream of computational fluid dynamics, for example, in gas dynamics, compressible and incompressible flows, turbomachinery, magneto-hydrodynamics, granular flows, semiconductor device simulation, etc., but have also prompted their application to a wide variety of problems for which they were not originally intended like, for example, second-order elliptic problems, and elasticity.

Some more stuff here, references, etc...

Chapter 1

Compressible flow and the Euler equations

1.1 Equations describing the flow

We consider a time interval $(0, T)$ and space domain $\Omega_t \subset \mathbb{R}^3$ occupied by a fluid at time t . By \mathcal{M} we denote the space-time domain in consideration:

$$\mathcal{M} = \{(\mathbf{x}, t); \mathbf{x} \in \Omega_t, t \in (0, T)\}. \quad (1.1)$$

Moreover we assume that \mathcal{M} is an open set.

1.1.1 Description of the flow

In Computational Fluid Dynamics there exist two classical approaches to the description of the flow, the *Lagrangian description* and the *Eulerian description*.

The idea of the Lagrangian description is to monitor each fluid particle along its pathline (i.e. the curve which the particle traverses in time). If we wanted to set a computational mesh using this description, it would mean to firmly connect nodes of the mesh with certain particles (i.e. the node and the particle would have to share their space coordinates) and move the mesh accordingly to the motion of the fluid as to preserve the node - corresponding particle connection at each time instant. The obvious drawback is the necessity to perform re-meshing operations very frequently, especially when dealing with large distortion of the fluid.

The Eulerian description focuses on fluid particles that move through fixed points within a computational domain. In other words, whereas in the Lagrangian description the particle was fixed and the point in space it was currently occupying was changing, now it is the point in space that holds still and the particle in consideration is changing and is always corresponding to the one that is currently occupying

the considered point in space. From this idea it follows that a computational mesh for this description would be fixed with respect to time.

It is not difficult to imagine that formulation of some basic mechanical principles could be easier for the moving particle, that is using the Lagrangian approach. However, the Eulerian description is used for the formulation of conservation laws as will be seen in the following subsections. Last, using this approach, large distortions of fluid domain can be handled with relative ease.

We shall proceed now with basics of the Lagrangian and the Eulerian descriptions and their relation. We shall present the equations describing the flow derived from conservation laws in their integral forms using the Eulerian approach.

Lagrangian description

We specify the particle in consideration using the mapping

$$\varphi(\mathbf{X}, t_0; t) \quad (1.2)$$

which determines the current (at time t) position $\mathbf{x} \in \Omega_t$ of the particle that occupies the point \mathbf{X} at time t_0 , i.e.

$$\mathbf{x} = \varphi(\mathbf{X}, t_0; t), \quad \mathbf{X} \in \Omega_{t_0}, \quad (1.3)$$

where we can omit the reference time t_0 and write

$$\mathbf{x} = \varphi(\mathbf{X}, t). \quad (1.4)$$

Customarily the components X_1, X_2, X_3 of the reference point \mathbf{X} are called the *Lagrangian coordinates* and the components x_1, x_2, x_3 of the point \mathbf{x} in the current configuration Ω_t are called the *Eulerian coordinates*. The *velocity* and *acceleration* of the particle given by the reference point \mathbf{X} are defined as

$$\hat{v}(\mathbf{X}, t) = \frac{\partial \varphi}{\partial t}(\mathbf{X}, t_0; t), \quad (1.5)$$

$$\hat{a}(\mathbf{X}, t) = \frac{\partial^2 \varphi}{\partial t^2}(\mathbf{X}, t_0; t), \quad (1.6)$$

provided the above derivatives exist.

Eulerian description

Using once again the mapping φ defined in (1.2), relation (1.4) and the Lagrangian definition of velocity (1.5), we can express the *velocity* of the fluid particle passing through the point \mathbf{x} at time t :

$$\mathbf{v}(\mathbf{x}, t) = \hat{v}(\mathbf{X}, t) = \frac{\partial \varphi}{\partial t}(\mathbf{X}, t), \quad (1.7)$$

where $\mathbf{x} = \boldsymbol{\varphi}(\mathbf{X}, t)$.

We shall demand the following regularity from the velocity function:

$$\mathbf{v} \in [\mathcal{C}^1(\mathcal{M})]^3. \quad (1.8)$$

We can pass to the Eulerian coordinates from the Lagrangian ones by solving the following initial value problem:

$$\frac{\partial \mathbf{x}}{\partial t} = \mathbf{v}(\mathbf{x}, t), \quad \mathbf{x}(t_0) = \mathbf{X}. \quad (1.9)$$

Under assumption (1.8), the problem (1.9) has exactly one maximal solution $\boldsymbol{\varphi}(\mathbf{X}, t_0, t)$ for each $(\mathbf{X}, t_0) \in \mathcal{M}$ defined for t from a certain subinterval of $(0, T)$. Moreover, in its domain of definition, the mapping $\boldsymbol{\varphi}$ has continuous first order derivatives with respect to X_1, X_2, X_3, t_0, t and continuous second order derivatives $\partial^2 \boldsymbol{\varphi} / \partial t \partial X_i, \partial^2 \boldsymbol{\varphi} / \partial t_0 \partial X_i, i = 1, 2, 3$. These statements result from theory of classical solutions of ordinary differential equations.

Under assumption (1.8), the *acceleration* of the particle passing through the point \mathbf{x} at time t can be expressed as

$$\mathbf{a}(\mathbf{x}, t) = \frac{\partial \mathbf{v}}{\partial t}(\mathbf{x}, t) + \sum_{i=1}^3 v_i(\mathbf{x}, t) \frac{\partial \mathbf{v}}{\partial x_i}(\mathbf{x}, t). \quad (1.10)$$

This, written in a short form reads

$$\mathbf{a} = \frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \text{grad}) \mathbf{v} = \frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla) \mathbf{v}, \quad (1.11)$$

where the differentiation represented by the symbol ∇ is with respect to the spatial variables x_1, x_2, x_3 .

1.1.2 The Transport theorem

We wish to study some physical quantity which is transported by fluid particles in our space-time domain \mathcal{M} . Let a function $F = F(\mathbf{x}, t) : \mathcal{M} \rightarrow \mathbb{R}$ represent some physical quantity in the Eulerian coordinates, and let us consider a system of fluid particles filling a bounded domain $\mathcal{V}(t) \subset \Omega_t$ at time t . By \mathcal{F} we denote the total amount of the quantity represented by the function F contained in $\mathcal{V}(t)$:

$$\mathcal{F}(t) = \int_{\mathcal{V}(t)} F(\mathbf{x}, t) d\mathbf{x}. \quad (1.12)$$

For the formulation of fundamental equations describing the flow we need to calculate the rate of change of the quantity \mathcal{F} bound on the system of particles considered. In other words we shall be interested in the derivative

$$\frac{d\mathcal{F}(t)}{dt} = \frac{d}{dt} \int_{\mathcal{V}(t)} F(\mathbf{x}, t) d\mathbf{x}. \quad (1.13)$$

In the following we shall suppose that (1.8) holds.

Lemma 1.1

Let $t_0 \in (0, T)$, $\mathcal{V}(t_0)$ be a bounded domain and let $\overline{\mathcal{V}(t_0)} \subset \Omega_{t_0}$. Then there exist an interval $(t_1, t_2) \subset (0, T)$, $t_0 \in (t_1, t_2)$ such that the following conditions are satisfied:

- a) The mapping ' $t \in (t_1, t_2), \mathbf{X} \in \mathcal{V}(t_0) \longrightarrow \mathbf{x} = \boldsymbol{\varphi}(\mathbf{X}, t_0; t) \in \mathcal{V}(t)$ ' has continuous first order derivatives with respect to t, X_1, X_2, X_3 and continuous second order derivatives $\partial^2 \boldsymbol{\varphi} / \partial t \partial X_i$, $i = 1, 2, 3$.
- b) The mapping ' $\mathbf{X} \in \mathcal{V}(t_0) \longrightarrow \mathbf{x} = \boldsymbol{\varphi}(\mathbf{X}, t_0; t) \in \mathcal{V}(t)$ ' is for all $t \in (t_1, t_2)$ a continuously differentiable one-to-one mapping of $\mathcal{V}(t_0)$ onto $\mathcal{V}(t)$ with continuous and bounded Jacobian $\mathcal{J}(\mathbf{X}, t)$ which satisfies the condition

$$\mathcal{J}(\mathbf{X}, t) > 0 \quad \forall \mathbf{X} \in \mathcal{V}(t_0), \quad \forall t \in (t_1, t_2).$$

- c) The inclusion

$$\left\{ (\mathbf{x}, t); t \in [t_1, t_2], \mathbf{x} \in \overline{\mathcal{V}(t)} \right\} \subset \mathcal{M}$$

holds and therefore the mapping \mathbf{v} has continuous and bounded first order derivatives on $\{(\mathbf{x}, t); t \in (t_1, t_2), \mathbf{x} \in \mathcal{V}(t)\}$ with respect to all variables.

- d) $\mathbf{v}(\boldsymbol{\varphi}(\mathbf{X}, t_0; t), t) = \frac{\partial \boldsymbol{\varphi}}{\partial t}(\mathbf{X}, t_0; t) \quad \forall \mathbf{X} \in \mathcal{V}(t_0), \quad \forall t \in (t_1, t_2).$

For proof, see [?].

Theorem 1.2 - The transport theorem

Let conditions from Lemma 1.1, a)-d) be satisfied and let the function $F = F(\mathbf{x}, t)$ have continuous and bounded first order derivatives on the set $\{(\mathbf{x}, t); t \in (t_1, t_2), \mathbf{x} \in \mathcal{V}(t)\}$. Let $\mathcal{F}(t) = \int_{\mathcal{V}(t)} F(\mathbf{x}, t) d\mathbf{x}$. Then for each $t \in (t_1, t_2)$ there exists a finite derivative

$$\begin{aligned} \frac{d\mathcal{F}}{dt}(t) &= \frac{d}{dt} \int_{\mathcal{V}(t)} F(\mathbf{x}, t) d\mathbf{x} \\ &= \int_{\mathcal{V}(t)} \left[\frac{\partial F}{\partial t}(\mathbf{x}, t) + \operatorname{div}(F\mathbf{v})(\mathbf{x}, t) \right] d\mathbf{x}. \end{aligned} \quad (1.14)$$

For proof, see [?].

1.1.3 The continuity equation

The density of fluid is a function

$$\rho : \mathcal{M} \longrightarrow (0, +\infty)$$

which allows us to determine the mass $m(\mathcal{V}; t)$ of the fluid contained in any sub-domain $\mathcal{V} \subset \Omega_t$:

$$m(\mathcal{V}; t) = \int_{\mathcal{V}} \rho(\mathbf{x}, t) d\mathbf{x}. \quad (1.15)$$

Assumptions 1.3

In what follows, let $\rho \in \mathcal{C}^1(\mathcal{M})$ and as before let $\mathbf{v} \in [\mathcal{C}^1(\mathcal{M})]^3$. We shall consider an arbitrary time instant $t_0 \in (0, T)$ and a moving piece of fluid formed by the same particles at each time instant and filling at time t_0 a bounded domain $\mathcal{V} \subset \overline{\mathcal{V}} \subset \Omega_{t_0}$ with a Lipschitz-continuous boundary $\partial\mathcal{V}$ called the *control volume* in the domain Ω_{t_0} . By $\mathcal{V}(t)$ we denote the domain occupied by this piece of fluid at time $t \in (t_1, t_2)$, where (t_1, t_2) is a sufficiently small interval containing t_0 with properties from Lemma 1.1.

Since the domain $\mathcal{V}(t)$ is formed by the same particles at each time instant, the *conservation of mass* can be formulated in the following way: *The mass of the piece of fluid represented by the domain $\mathcal{V}(t)$ does not depend on time t .* This means that

$$\frac{dm(\mathcal{V}(t); t)}{dt} = 0, \quad t \in (t_1, t_2), \quad (1.16)$$

where with respect to (1.15) we have

$$m(\mathcal{V}(t); t) = \int_{\mathcal{V}(t)} \rho(\mathbf{x}, t) d\mathbf{x}. \quad (1.17)$$

From Theorem 1.2 for a function $F := \rho$ we get the identity

$$\int_{\mathcal{V}(t)} \left[\frac{\partial \rho}{\partial t}(\mathbf{x}, t) + \operatorname{div}(\rho \mathbf{v})(\mathbf{x}, t) \right] d\mathbf{x} = 0, \quad t \in (t_1, t_2). \quad (1.18)$$

If we substitute $t := t_0$ and take into account that $\mathcal{V}(t_0) = \mathcal{V}$, we conclude that

$$\int_{\mathcal{V}} \left[\frac{\partial \rho}{\partial t}(\mathbf{x}, t_0) + \operatorname{div}(\rho \mathbf{v})(\mathbf{x}, t_0) \right] d\mathbf{x} = 0 \quad (1.19)$$

for an arbitrary $t_0 \in (0, T)$ and an arbitrary control volume $\mathcal{V} \subset \Omega_{t_0}$. We use the following Lemma in order to derive the differential form of the law of conservation of mass:

Lemma 1.4

Let $\Omega \subset \mathbb{R}^N$ be an open set and let $f \in \mathcal{C}(\Omega)$. Then the following holds:

$f \equiv 0$ in Ω if and only if $\int_{\mathcal{V}} f(\mathbf{x}) d\mathbf{x} = 0$ for any bounded open set $\mathcal{V} \subset \overline{\mathcal{V}} \subset \Omega$.

Now we use Lemma 1.4 and obtain the differential form of the law of conservation of mass called the *continuity equation*:

$$\frac{\partial \rho}{\partial t}(\mathbf{x}, t) + \operatorname{div}(\rho(\mathbf{x}, t) \mathbf{v}(\mathbf{x}, t)) = 0, \quad \mathbf{x} \in \Omega_t, \quad t \in (0, T). \quad (1.20)$$

1.1.4 The equations of motion

We proceed by deriving basic dynamical equations describing flow motion from the *law of conservation of momentum* which can be formulated in this way:

The rate of change of the total momentum of a piece of fluid formed by the same particles at each time and occupying the domain $\mathcal{V}(t)$ at time instant t is equal to the force acting on $\mathcal{V}(t)$.

Let assumptions 1.3 be satisfied. The total momentum of particles contained in $\mathcal{V}(t)$ is given by

$$\mathcal{H}(\mathcal{V}(t)) = \int_{\mathcal{V}(t)} \rho(\mathbf{x}, t) \mathbf{v}(\mathbf{x}, t) d\mathbf{x}. \quad (1.21)$$

Moreover, denoting by $\mathcal{F}(\mathcal{V}(t))$ the force acting on the volume \mathcal{V} , the law of conservation of momentum reads

$$\frac{d\mathcal{H}(\mathcal{V}(t))}{dt} = \mathcal{F}(\mathcal{V}(t)), \quad t \in (t_1, t_2). \quad (1.22)$$

Using Theorem 1.2 for functions $F := \rho v_i$, $i = 1, 2, 3$, we get

$$\int_{\mathcal{V}(t)} \left[\frac{\partial}{\partial t} (\rho(\mathbf{x}, t) v_i(\mathbf{x}, t)) + \operatorname{div}(\rho(\mathbf{x}, t) v_i(\mathbf{x}, t) \mathbf{v}(\mathbf{x}, t)) \right] d\mathbf{x} = \mathcal{F}_i(\mathcal{V}(t)), \quad (1.23)$$

$$i = 1, 2, 3, \quad t \in (t_1, t_2).$$

Taking into account that $t_0 \in (0, T)$ is an arbitrary time instant and $\mathcal{V}_{t_0} = \mathcal{V} \subset \overline{\mathcal{V}} \subset \Omega_{t_0}$, where \mathcal{V} is an arbitrary control volume, we get the law of conservation of momentum in the form where we write t instead of t_0 :

$$\int_{\mathcal{V}} \left[\frac{\partial}{\partial t} (\rho(\mathbf{x}, t) v_i(\mathbf{x}, t)) + \operatorname{div}(\rho(\mathbf{x}, t) v_i(\mathbf{x}, t) \mathbf{v}(\mathbf{x}, t)) \right] d\mathbf{x} = \mathcal{F}_i(\mathcal{V}; t), \quad (1.24)$$

$i = 1, 2, 3$, for an arbitrary $t \in (0, T)$ and an arbitrary control volume $\mathcal{V} \subset \overline{\mathcal{V}} \subset \Omega_t$.

According to [?], the components $\mathcal{F}_i(\mathcal{V}; t)$, $i = 1, 2, 3$, of the vector $\mathcal{F}(\mathcal{V}; t)$ can be expressed as

$$\mathcal{F}_i(\mathcal{V}; t) = \int_{\mathcal{V}} \rho(\mathbf{x}, t) f_i(\mathbf{x}, t) dx + \int_{\partial\mathcal{V}} \sum_{j=1}^3 \tau_{ji}(\mathbf{x}, t) n_j(\mathbf{x}) dS, \quad i = 1, 2, 3, \quad (1.25)$$

assuming that $\tau_{ij} \in \mathcal{C}^1(\mathcal{M})$ and $f_i \in \mathcal{C}(\mathcal{M})$, $(i, j = 1, 2, 3)$. Here τ_{ji} are components of the *stress tensor* \mathcal{T} and f_i are components of the *density of the volume force* \mathbf{f} . Substituting this into (1.24), we get

$$\int_{\mathcal{V}} \left[\frac{\partial}{\partial t} (\rho(\mathbf{x}, t) v_i(\mathbf{x}, t)) + \operatorname{div} (\rho(\mathbf{x}, t) v_i(\mathbf{x}, t) \mathbf{v}(\mathbf{x}, t)) \right] dx = \quad (1.26)$$

$$\int_{\mathcal{V}} \rho(\mathbf{x}, t) f_i(\mathbf{x}, t) dx + \int_{\partial\mathcal{V}} \sum_{j=1}^3 \tau_{ji}(\mathbf{x}, t) n_j(\mathbf{x}) dS, \quad i = 1, 2, 3, \quad (1.27)$$

for each $t \in (0, T)$ and an arbitrary control volume \mathcal{V} in Ω_t . Moreover, applying Green's theorem and Lemma 1.4, we obtain the desired *equation of motion of a general fluid in the differential conservative form*

$$\frac{\partial}{\partial t} (\rho v_i) + \operatorname{div} (\rho v_i \mathbf{v}) = \rho f_i + \sum_{j=1}^3 \frac{\partial \tau_{ji}}{\partial x_j}, \quad i = 1, 2, 3. \quad (1.28)$$

This can be written as

$$\frac{\partial}{\partial t} (\rho \mathbf{v}) + \operatorname{div} (\rho \mathbf{v} \otimes \mathbf{v}) = \rho \mathbf{f} + \operatorname{div} \mathcal{T}, \quad (1.29)$$

where \otimes denotes the *tensor product*:

$$\mathbf{a} \otimes \mathbf{b} = \begin{pmatrix} a_1 b_1 & a_1 b_2 & a_1 b_3 \\ a_2 b_1 & a_2 b_2 & a_2 b_3 \\ a_3 b_1 & a_3 b_2 & a_3 b_3 \end{pmatrix},$$

and $\operatorname{div} (\mathbf{a} \otimes \mathbf{b})$ is a vector quantity:

$$\operatorname{div} (\mathbf{a} \otimes \mathbf{b}) = \left(\sum_{i=1}^3 \frac{\partial}{\partial x_i} a_i b_1, \sum_{i=1}^3 \frac{\partial}{\partial x_i} a_i b_2, \sum_{i=1}^3 \frac{\partial}{\partial x_i} a_i b_3 \right)^T.$$

1.1.5 The Navier-Stokes equations

The relation between the stress tensor and other quantities describing fluid flow, the velocity and its derivatives in particular, represent the so-called *rheological*

equations of the fluid. For the derivation of the Navier-Stokes equations we shall use

$$\mathcal{T} = (-p + \lambda \operatorname{div} \mathbf{v}) \mathbb{I} + 2\mu \mathbb{D}(\mathbf{v}), \quad (1.30)$$

where \mathbb{D} is the deformation velocity tensor:

$$\mathbb{D} = \mathbb{D}(\mathbf{v}) = (d_{ij})_{i,j=1}^3, \quad d_{ij} = \frac{1}{2} \left(\frac{\partial v_i}{\partial x_j} + \frac{\partial v_j}{\partial x_i} \right), \quad (1.31)$$

λ, μ are constants or scalar functions of thermodynamical quantities, λ and μ are called the *first* and the *second viscosity coefficient* respectively. For the assumptions under which we can write (1.30) see [?]. Although viscosity coefficients can be functions of thermodynamical quantities (most important of which is θ , the absolute temperature) we shall treat them as if they were constants. Let assumptions 1.3 be satisfied and let us assume that

$$\frac{\partial \mathbf{v}}{\partial t} \in [\mathcal{C}(\mathcal{M})]^3, \quad \frac{\partial^2 \mathbf{v}}{\partial x_i \partial x_j} \in [\mathcal{C}(\mathcal{M})]^3 \quad (i, j = 1, 2, 3). \quad (1.32)$$

Now let us substitute relation (1.30) into the general equations of motion (1.29) with the assumption of constant viscosity coefficients and assumptions (1.32). We come to the Navier-Stokes equations in the form

$$\frac{\partial(\rho \mathbf{v})}{\partial t} + \operatorname{div}(\rho \mathbf{v} \otimes \mathbf{v}) = \rho \mathbf{f} - \nabla p + \mu \Delta \mathbf{v} + (\mu + \lambda) \nabla \operatorname{div} \mathbf{v}. \quad (1.33)$$

For details see [?].

1.1.6 The energy equation

Now let us derive a differential equation equivalent to the *law of conservation of energy*. As in the preceding subsections, we consider a piece of fluid represented by a control volume $\mathcal{V}(t)$ satisfying assumptions 1.3. The law of conservation of energy can be formulated as follows:

The rate of change of the total energy of the fluid particles, occupying the domain $\mathcal{V}(t)$ at time t , is equal to the sum of powers of the volume force acting on the volume $\mathcal{V}(t)$ and the surface force acting on the surface $\partial \mathcal{V}(t)$, and of the amount of heat transmitted to $\mathcal{V}(t)$.

By $\mathcal{E}(\mathcal{V}(t))$ let us denote the total energy of the fluid particles contained in the domain $\mathcal{V}(t)$ and by $\mathcal{Q}(\mathcal{V}(t))$ the amount of heat transmitted to $\mathcal{V}(t)$ at time t . Taking into account the character of volume and surface forces involved, we get the identity representing the law of conservation of energy:

$$\begin{aligned} \frac{d}{dt} \mathcal{E}(\mathcal{V}(t)) &= \int_{\mathcal{V}(t)} \rho(\mathbf{x}, t) \mathbf{f}(\mathbf{x}, t) \cdot \mathbf{v}(\mathbf{x}, t) d\mathbf{x} \\ &+ \int_{\partial \mathcal{V}(t)} \sum_{i,j=1}^3 \tau_{ji}(\mathbf{x}, t) n_j(\mathbf{x}) v_i(\mathbf{x}, t) dS + \mathcal{Q}(\mathcal{V}(t)). \end{aligned} \quad (1.34)$$

Following relations hold:

$$\begin{aligned}
a) \quad \mathcal{E}(\mathcal{V}(t)) &= \int_{\mathcal{V}(t)} E(\mathbf{x}, t) d\mathbf{x}, \\
b) \quad E &= \rho \left(e + \frac{|\mathbf{v}|^2}{2} \right), \\
c) \quad \mathcal{Q}(\mathcal{V}(t)) &= \int_{\mathcal{V}(t)} \rho(\mathbf{x}, t) q(\mathbf{x}, t) d\mathbf{x} - \int_{\partial\mathcal{V}(t)} \boldsymbol{\phi}_q(\mathbf{x}, t) \cdot \mathbf{n}(\mathbf{x}) dS.
\end{aligned} \tag{1.35}$$

Here E is the total energy, e is the density of the specific internal energy (related to the unit mass) associated with molecular and atomic behavior, $|\mathbf{v}|^2/2$ is the density of the kinetic energy, q represents the density of heat sources (again related to the unit mass) and $\boldsymbol{\phi}_q$ is the heat flux.

Let assumptions 1.3 hold and further let $\tau_{ij}, (\boldsymbol{\phi}_q)_i \in \mathcal{C}^1(\mathcal{M})$ and $f_i, q \in \mathcal{C}(\mathcal{M})$ ($i, j = 1, 2, 3$). Using this, relations (1.35) a)-c), Theorem 1.2, Green's theorem and Lemma 1.4, we derive from (1.34) the differential *energy equation*, where we take advantage of (1.30):

$$\frac{\partial E}{\partial t} + \operatorname{div}(E\mathbf{v}) = \rho \mathbf{f} \cdot \mathbf{v} - \operatorname{div}(p\mathbf{v}) + \operatorname{div}(\lambda \mathbf{v} \operatorname{div} \mathbf{v}) + \operatorname{div}(2\mu \mathbb{D}(\mathbf{v}) \mathbf{v}) + \rho q - \operatorname{div} \boldsymbol{\phi}_q. \tag{1.36}$$

For details see [?].

1.1.7 Thermodynamical relations

In order to complete the equations describing the flow, some other relations shall be added. The system now contains seven unknown quantities: $v_1, v_2, v_3, \rho, e, \theta, p$, but only 5 equations (scalar continuity equation, vector Navier-Stokes equations and scalar energy equation), i.e. $(1 + 3 + 1) = 5$. From this we see, that additional two equations should be included.

Basic Thermodynamical Quantities

The absolute temperature θ , the density ρ and the pressure p are called the *state variables*. All these quantities are positive scalar functions. We consider only the so-called *perfect gas* or *ideal gas* whose state variables satisfy the following *equation of state*

$$p = R\theta\rho, \tag{1.37}$$

where R is the *gas constant*, which is defined as

$$R = c_p - c_v. \tag{1.38}$$

Here c_p denotes the *specific heat at constant pressure*, i.e. the ratio of the increment of the amount of heat related to the unit mass, to the increment of temperature at constant pressure. Analogously c_v denotes the *specific heat at constant volume*. Experiments show that $c_p > c_v$, so that $R > 0$, and that c_p and c_v can be treated like constants for a relatively large range of temperature. We set $\gamma = c_p/c_v$ which is the so-called *Poisson adiabatic constant*. The internal energy related to the unit mass is defined by

$$e = c_v \theta, \quad (1.39)$$

which explains the meaning of the internal energy: it is the amount of heat it would have to be transmitted out of the fluid so that its temperature would reach (absolute) zero, volume being kept constant during the whole process.

With respect to the above relations, we can express the internal energy as

$$e = c_p \theta - R \theta. \quad (1.40)$$

The complete system of equations describing the flow

The complete system now reads

$$\frac{\partial \rho}{\partial t} + \operatorname{div}(\rho \mathbf{v}) = 0, \quad (1.41)$$

$$\frac{\partial(\rho \mathbf{v})}{\partial t} + \operatorname{div}(\rho \mathbf{v} \otimes \mathbf{v}) = \rho \mathbf{f} - \nabla p + \mu \Delta \mathbf{v} + (\mu + \lambda) \nabla \operatorname{div} \mathbf{v}, \quad (1.42)$$

$$\begin{aligned} \frac{\partial E}{\partial t} + \operatorname{div}(E \mathbf{v}) &= \rho \mathbf{f} \cdot \mathbf{v} - \operatorname{div}(p \mathbf{v}) + \operatorname{div}(\lambda \mathbf{v} \operatorname{div} \mathbf{v}) + \\ &+ \operatorname{div}(2\mu \mathbb{D}(\mathbf{v})) + \rho q - \operatorname{div} \boldsymbol{\phi}_q, \end{aligned} \quad (1.43)$$

$$p = (\gamma - 1) (E - \rho |\mathbf{v}|^2 / 2), \quad (1.44)$$

$$\theta = (E/\rho - |\mathbf{v}|^2 / 2) / c_v. \quad (1.45)$$

This system is simply called the *compressible Navier-Stokes equations* for a heat-conductive perfect gas. Equations (1.44) and (1.45) follow from (1.37) - (1.40) and (1.35).

1.1.8 Entropy and the second law of thermodynamics

Entropy One of the important thermodynamical quantities is the entropy S , defined by the relation

$$\theta dS = de + p dV, \quad (1.46)$$

where $V = 1/\rho$ is the so-called specific volume. This identity is derived in thermodynamics under the assumption that the internal energy is a function of S and V , $e = e(S, V)$, which explains the meaning of the differentials in (1.46).

Theorem 1.3

For a perfect gas we have

$$S = c_v \ln \frac{p/p_0}{(\rho/\rho_0)^\gamma} + \text{const} \quad (1.47)$$

$$= c_v \ln \frac{\theta/\theta_0}{(\rho/\rho_0)^{\gamma-1}} + \text{const}, \quad (1.48)$$

where p_0 and ρ_0 are fixed (reference) values of pressure and density and $\theta_0 = p_0/(R\rho_0)$.

Proof

Using (1.39) and the relation $V = 1/\rho$, we can write (1.46) in the form

$$\theta dS = c_v d\theta - \frac{p d\rho}{\rho^2}. \quad (1.49)$$

From this and (1.37) - (1.39) we get

$$dS = c_v \frac{d\theta}{\theta} - \frac{p}{\rho\theta} \frac{d\rho}{\rho} = c_v \frac{d(p/\rho)}{(p/\rho)} - R \frac{d\rho}{\rho} = c_v d \ln \frac{p/p_0}{(\rho/\rho_0)^\gamma} = c_v d \ln \frac{\theta/\theta_0}{(\rho/\rho_0)^{\gamma-1}}, \quad (1.50)$$

which immediately yields (1.47).

The second law of thermodynamics In the irreversible processes, equality (1.47) does not hold in general and is replaced by the inequality

$$dS \geq \frac{\delta Q}{\theta} \quad (1.51)$$

called the *second law of thermodynamics*. For a system of fluid particles occupying a domain $\mathcal{V}(t)$ at time t we postulate the second law of thermodynamics mathematically in the form

$$\frac{d}{dt} \int_{\mathcal{V}(t)} \rho(\mathbf{x}, t) S(\mathbf{x}, t) d\mathbf{x} \geq \int_{\mathcal{V}(t)} \frac{\rho(\mathbf{x}, t) q(\mathbf{x}, t)}{\theta(\mathbf{x}, t)} d\mathbf{x} - \int_{\delta\mathcal{V}(t)} \frac{\phi_q(\mathbf{x}, t) \cdot \mathbf{n}(\mathbf{x})}{\theta(\mathbf{x}, t)} dS. \quad (1.52)$$

The left-hand side of (1.52) represents the rate of change of the entropy contained in the volume $\mathcal{V}(t)$, and the first and second integral on the right-hand side are called the *entropy production* and the *entropy flux*. Let $\rho, \theta, v_i, \phi_{qi} \in \mathcal{C}^1(\mathcal{M})$; $q, f_i \in \mathcal{C}(\mathcal{M})$, $i = 1, 2, 3$. By virtue of the transport Theorem 1.2 and the continuity equation (1.20), from (1.52) we obtain the inequality

$$\rho \frac{\partial(\rho S)}{\partial t} + \text{div}(\rho S \mathbf{v}) \geq \frac{\rho q}{\theta} - \text{div} \left(\frac{\phi_q}{\theta} \right). \quad (1.53)$$

A more general model of flow is obtained in thermodynamics under the assumption that the pressure is a function of the density and entropy: $p = p(\rho, S)$, where p is a continuously differentiable function and $\partial p / \partial \rho > 0$. Let us introduce the quantity

$$a = \sqrt{\frac{\partial p}{\partial \rho}} \quad (1.54)$$

which has the dimension ms^{-1} of velocity and is called the *speed of sound*. Another important characteristic of the flow is so-called *Mach number*, which is defined as

$$M = \frac{|v|}{a} \quad , \quad (1.55)$$

where v is the flow velocity.

1.2 Euler equations and their properties

If we set $\mu = \lambda = k = 0$, we obtain the model of inviscid compressible flow, described by the continuity equation, the Euler equations, the energy equation and thermodynamical relations. Since gases are light, usually it is possible to neglect the effect of the volume force. Neglecting heat sources too (assuming *adiabatic flow*), we get the system for the perfect inviscid gas in the following form:

$$\frac{\partial \rho}{\partial t} + \operatorname{div}(\rho \mathbf{v}) = 0, \quad (1.56)$$

$$\frac{\partial(\rho \mathbf{v})}{\partial t} + \operatorname{div}(\rho \mathbf{v} \otimes \mathbf{v}) + \nabla p = 0, \quad (1.57)$$

$$\frac{\partial E}{\partial t} + \operatorname{div}((E + p) \mathbf{v}) = 0, \quad (1.58)$$

$$p = (\gamma - 1) (E - \rho |\mathbf{v}|^2 / 2). \quad (1.59)$$

This system is simply called the *compressible Euler equations*.

In the following we shall be concerned with the 2-dimensional case, which will be used in numerical examples. All results apply for the 3-dimensional case as well.

1.2.1 Dimensionless Euler equations

We set 3 similarity constants l_r, v_r , and ρ_r , which are, in order: characteristic length, characteristic velocity, and characteristic density. Now we multiply the Euler equations by these constants in the following way:

$$\left[\frac{\partial \rho}{\partial t} + \operatorname{div}(\rho \mathbf{v}) \right] \frac{l_r}{\rho_r v_r} = 0, \quad (1.60)$$

$$\left[\frac{\partial(\rho \mathbf{v})}{\partial t} + \operatorname{div}(\rho \mathbf{v} \otimes \mathbf{v}) + \nabla p \right] \frac{l_r}{\rho_r v_r^2} = 0, \quad (1.61)$$

$$\left[\frac{\partial E}{\partial t} + \operatorname{div}((E + p) \mathbf{v}) \right] \frac{l_r}{\rho_r v_r^3} = 0. \quad (1.62)$$

Then we get the dimensionless Euler equations in the form

$$\frac{\partial \tilde{\rho}}{\partial \tilde{t}} + \tilde{\nabla} \cdot (\tilde{\rho} \tilde{\mathbf{v}}) = 0, \quad (1.63)$$

$$\frac{\partial(\tilde{\rho} \tilde{\mathbf{v}})}{\partial \tilde{t}} + \tilde{\nabla} \cdot (\tilde{\rho} \tilde{\mathbf{v}} \otimes \tilde{\mathbf{v}}) + \tilde{\nabla} \tilde{p} = 0, \quad (1.64)$$

$$\frac{\partial \tilde{E}}{\partial \tilde{t}} + \tilde{\nabla} \cdot ((\tilde{E} + \tilde{p}) \tilde{\mathbf{v}}) = 0, \quad (1.65)$$

where

$$t_r = \frac{l_r}{v_r}, \quad (1.66)$$

$$\tilde{t} = \frac{t}{t_r}, \quad (1.67)$$

$$\tilde{\rho} = \frac{\rho}{\rho_r}, \quad (1.68)$$

$$\tilde{\mathbf{v}} = \frac{\mathbf{v}}{v_r}, \quad (1.69)$$

$$\tilde{\nabla} = l_r \nabla, \quad (1.70)$$

$$\tilde{E} = \frac{E}{\rho_r u_r^2}, \quad (1.71)$$

$$\tilde{p} = \frac{p}{\rho_r u_r^2}. \quad (1.72)$$

As we see, the dimensionless Euler equations have the same form as the original equations, only the relations are rescaled using the above relations. Therefore, in the sequel, we shall omit the symbol " " over the considered quantities.

1.2.2 Conservative form of the Euler equations

The system of governing equations can be written in the so-called *conservative form*.

$$\frac{\partial \mathbf{w}}{\partial t} + \frac{\partial \mathbf{f}_1}{\partial x_1} + \frac{\partial \mathbf{f}_2}{\partial x_2} = 0, \quad (1.73)$$

where

$$\mathbf{w} = \begin{pmatrix} \varrho \\ \rho v_1 \\ \rho v_2 \\ E \end{pmatrix} = \begin{pmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{pmatrix}, \quad (1.74)$$

$$\mathbf{f}_1 = \begin{pmatrix} \rho v_1 \\ \rho v_1^2 + p \\ \rho v_1 v_2 \\ v_1(E + p) \end{pmatrix} = \begin{pmatrix} w_2 \\ \frac{w_2^2}{w_1} + p \\ \frac{w_2 w_3}{w_1} \\ \frac{w_2}{w_1}(w_4 + p) \end{pmatrix}, \quad (1.75)$$

$$\mathbf{f}_2 = \begin{pmatrix} \rho v_2 \\ \rho v_2 v_1 \\ \rho v_2^2 + p \\ v_2(E + p) \end{pmatrix} = \begin{pmatrix} w_2 \\ \frac{w_3 w_2}{w_1} \\ \frac{w_3^2}{w_1} + p \\ \frac{w_3}{w_1}(w_4 + p) \end{pmatrix}, \quad (1.76)$$

$$p = (\gamma - 1) \left(w_4 - \frac{w_2^2 + w_3^2}{2w_1} \right). \quad (1.77)$$

Usually, \mathbf{f}_1 and \mathbf{f}_2 , are called *inviscid Euler fluxes*. The unknowns ρ, v_1, v_2, p in the original system of equations are called *primitive variables*, whereas $w_1 = \rho, w_2 = \rho v_1, w_3 = \rho v_2, w_4 = E$ are called *conservative variables*.

The domain of definition of the vector-valued functions $\mathbf{f}_1, \mathbf{f}_2$ is the open set $D \subset \mathbb{R}^2$ where the corresponding density and pressure are positive:

$$D = \left\{ \mathbf{w} \in \mathbb{R}^4; w_1 = \rho > 0, w_4 - \frac{w_2^2 + w_3^2}{2w_1} = \frac{p}{\gamma - 1} > 0 \right\}. \quad (1.78)$$

One can see that the inviscid Euler fluxes are continuously differentiable in their domain of definition. Differentiation in (1.73) and the chain rule lead to a first order system of quasilinear partial differential equations

$$\frac{\partial \mathbf{w}}{\partial t} + \mathbb{A}_1(\mathbf{w}) \frac{\partial \mathbf{w}}{\partial x_1} + \mathbb{A}_2(\mathbf{w}) \frac{\partial \mathbf{w}}{\partial x_2} = 0, \quad (1.79)$$

where $\mathbb{A}_s(\mathbf{w}), s = 1, 2$ are $m \times m$ matrices called *flux Jacobians* defined as the Jacobi matrices of the mappings $\mathbf{f}_s, s = 1, 2$, at a point $\mathbf{w} \in D$:

$$\mathbb{A}_s(\mathbf{w}) = \frac{D\mathbf{f}_s(\mathbf{w})}{D\mathbf{w}} = \left(\frac{\partial f_{si}(\mathbf{w})}{\partial w_j} \right)_{i,j=1}^2. \quad (1.80)$$

For $\mathbf{w} \in D$ and $\mathbf{n} = (n_1, n_2)^T \in \mathbb{R}^2$ we put

$$\mathcal{P}(\mathbf{w}, \mathbf{n}) = \sum_{s=1}^2 \mathbf{f}_s(\mathbf{w}) n_s, \quad (1.81)$$

which is the flux of the quantity \mathbf{w} in the direction \mathbf{n} . The Jacobi matrix $D\mathcal{P}(\mathbf{w}, \mathbf{n})/D\mathbf{w}$ can be expressed in the form

$$\frac{D\mathcal{P}(\mathbf{w}, \mathbf{n})}{D\mathbf{w}} = \mathbb{P}(\mathbf{w}, \mathbf{n}) = \sum_{s=1}^2 \mathbb{A}_s(\mathbf{w}) n_s. \quad (1.82)$$

Lemma 1.5

The vector-valued functions $\mathbf{f}_s, s = 1, 2$, defined by (1.75) - (1.76), are first order homogenous mappings:

$$\mathbf{f}_s(\alpha \mathbf{w}) = \alpha \mathbf{f}_s(\mathbf{w}), \quad \alpha > 0. \quad (1.83)$$

Moreover, we have

$$\mathbf{f}_s(\mathbf{w}) = \mathbb{A}_s(\mathbf{w}) \mathbf{w}. \quad (1.84)$$

Similarly,

$$\mathcal{P}(\alpha \mathbf{w}, \mathbf{n}) = \alpha \mathcal{P}(\mathbf{w}, \mathbf{n}), \quad \alpha > 0, \quad (1.85)$$

$$\mathcal{P}(\mathbf{w}, \mathbf{n}) = \mathbb{P}(\mathbf{w}, \mathbf{n}) \mathbf{w}. \quad (1.86)$$

Proof: Relation (1.83) immediately follows from (1.75) - (1.76). Since $\mathbf{f}_s \in \mathcal{C}^1(D)^m$, the expression $(D\mathbf{f}_s(\mathbf{w})/D\mathbf{w}) \mathbf{w} = \mathbb{A}_s(\mathbf{w}) \mathbf{w}$ is the derivative of \mathbf{f}_s in the direction \mathbf{w} at the point \mathbf{w} . By the definition of the derivative and (1.83),

$$\mathbb{A}_s(\mathbf{w}) \mathbf{w} = \lim_{\alpha \rightarrow 0} \frac{\mathbf{f}_s(\mathbf{w} + \alpha \mathbf{w}) - \mathbf{f}_s(\mathbf{w})}{\alpha} \quad (1.87)$$

$$= \lim_{\alpha \rightarrow 0} \frac{(1 + \alpha) \mathbf{f}_s(\mathbf{w}) - \mathbf{f}_s(\mathbf{w})}{\alpha} = \mathbf{f}_s(\mathbf{w}). \quad (1.88)$$

The rest of the Lemma follows from the definitions of \mathcal{P} and \mathbb{P} and the above results. The flux Jacobians have the following form:

$$\mathbb{A}_1(\mathbf{w}) = \frac{D\mathbf{f}_1}{D\mathbf{w}}$$

$$= \begin{pmatrix} 0 & 1 & 0 & 0 \\ -\frac{w_2^2}{w_1^2} + \frac{R}{c_v} \frac{w_2^2 + w_3^2}{2w_1^2} & \frac{2w_2}{w_1} - \frac{R}{c_v} \frac{w_2}{w_1} & -\frac{R}{c_v} \frac{w_3}{w_1} & \frac{R}{c_v} \\ -\frac{w_2 w_3}{w_1^2} & \frac{w_3}{w_1} & \frac{w_2}{w_1} & 0 \\ -\frac{w_2 w_4}{w_1^2} - \frac{w_2}{w_1} \frac{R}{c_v} \left(w_4 - \frac{w_2^2 + w_3^2}{2w_1} \right) + \frac{w_2}{w_1} \frac{R}{c_v} \frac{w_2^2 + w_3^2}{2w_1^2} & \frac{w_4}{w_1} + \frac{1}{w_1} \frac{R}{c_v} \left(w_4 - \frac{w_2^2 + w_3^2}{2w_1} \right) - \frac{R}{c_v} \frac{w_2^2}{w_1^2} & -\frac{R}{c_v} \frac{w_2 w_3}{w_1^2} & \frac{w_2}{w_1} + \frac{R}{c_v} \frac{w_2}{w_1} \end{pmatrix},$$

$$\mathbb{A}_2(\mathbf{w}) = \frac{D\mathbf{f}_2}{D\mathbf{w}}$$

$$= \begin{pmatrix} 0 & 0 & 1 & 0 \\ -\frac{w_3 w_2}{w_1^2} & \frac{w_3}{w_1} & \frac{w_2}{w_1} & 0 \\ -\frac{w_3^2}{w_1^2} + \frac{R}{c_v} \frac{w_2^2 + w_3^2}{2w_1^2} & -\frac{R}{c_v} \frac{w_2}{w_1} & \frac{2w_3}{w_1} - \frac{R}{c_v} \frac{w_3}{w_1} & \frac{R}{c_v} \\ -\frac{w_3 w_4}{w_1^2} - \frac{w_3}{w_1} \frac{R}{c_v} \left(w_4 - \frac{w_2^2 + w_3^2}{2w_1} \right) + \frac{w_3}{w_1} \frac{R}{c_v} \frac{w_2^2 + w_3^2}{2w_1^2} & -\frac{R}{c_v} \frac{w_3 w_2}{w_1^2} & \frac{w_4}{w_1} + \frac{1}{w_1} \frac{R}{c_v} \left(w_4 - \frac{w_2^2 + w_3^2}{2w_1} \right) - \frac{R}{c_v} \frac{w_3^2}{w_1^2} & \frac{w_3}{w_1} + \frac{R}{c_v} \frac{w_3}{w_1} \end{pmatrix}.$$

1.2.3 Rotational invariance

The rotational invariance of the Euler equations is represented by the relations

$$\mathcal{P}(\mathbf{w}, \mathbf{n}) = \sum_{s=1}^2 \mathbf{f}_s(\mathbf{w}) n_s = \mathbb{Q}^{-1}(\mathbf{n}) \mathbf{f}_1(\mathbb{Q}(\mathbf{n}) \mathbf{w}), \quad (1.89)$$

$$\mathbb{P}(\mathbf{w}, \mathbf{n}) = \sum_{s=1}^2 \mathbb{A}_s(\mathbf{w}) n_s = \mathbb{Q}^{-1}(\mathbf{n}) \mathbb{A}_1(\mathbb{Q}(\mathbf{n}) \mathbf{w}) \mathbb{Q}(\mathbf{n}), \quad (1.90)$$

$$\mathbf{n} = (n_1, n_2) \in \mathbb{R}^2, |\mathbf{n}| = 1, \mathbf{w} \in D, \quad (1.91)$$

where

$$\mathbb{Q}(\mathbf{n}) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & n_1 & n_2 & 0 \\ 0 & -n_2 & n_1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}. \quad (1.92)$$

1.2.4 Diagonalization of the Jacobi matrix

We have

$$\begin{aligned} \mathbf{f}_1(\mathbf{w}) &= (w_1, w_2^2/w_1 + (\gamma - 1) [w_4 - (w_2^2 + w_3^2) / (2w_1)], \\ &\quad w_2 w_3 / w_1, w_2 [\gamma w_4 - (\gamma - 1) (w_2^2 + w_3^2) / (2w_1)] / w_2)^T, \end{aligned} \quad (1.93)$$

and, with the notation $\mathbf{v} = (u, v)$,

$$\mathbb{A}_1(\mathbf{w}) = \begin{pmatrix} 0 & 1 & 0 & 0 \\ \frac{\gamma-1}{2} |\mathbf{v}|^2 - u^2 & (3-\gamma)u & (1-\gamma)v & 1-\gamma \\ -uv & v & u & 0 \\ u \left((\gamma-1) |\mathbf{v}|^2 - \gamma \frac{E}{\rho} \right) & \gamma \frac{E}{\rho} - (\gamma-1)u^2 - \frac{\gamma-1}{2} |\mathbf{v}|^2 & (1-\gamma)uv & \gamma u \end{pmatrix}. \quad (1.94)$$

The matrix $\mathbb{A}_1(\mathbf{w})$ has eigenvalues

$$\tilde{\lambda}_1(\mathbf{w}) = u - a, \tilde{\lambda}_2(\mathbf{w}) = \tilde{\lambda}_3(\mathbf{w}) = u, \tilde{\lambda}_4(\mathbf{w}) = u + a \quad (1.95)$$

and the corresponding eigenvectors are

$$\mathbf{r}_1(\mathbf{w}) = (1, u - a, v, |\mathbf{v}|^2 / 2 + a^2 (\gamma - 1) - ua)^T, \quad (1.96)$$

$$\mathbf{r}_2(\mathbf{w}) = (1, u, v, |\mathbf{v}|^2 / 2)^T, \quad (1.97)$$

$$\mathbf{r}_3(\mathbf{w}) = (1, u, v - a, |\mathbf{v}|^2 / 2 - va)^T, \quad (1.98)$$

$$\mathbf{r}_4(\mathbf{w}) = (1, u + a, v, |\mathbf{v}|^2 / 2 + a^2 (\gamma - 1) + ua)^T. \quad (1.99)$$

These calculations are possible to be verified. As the eigenvectors $\mathbf{r}_s(\mathbf{w})$, $s = 1, \dots, 4$ are linearly independent for each $\mathbf{w} \in D$, we can write

$$\tilde{\mathbb{T}}^{-1}(\mathbf{w}) \mathbb{A}_1(\mathbf{w}) \tilde{\mathbb{T}}(\mathbf{w}) = \tilde{\Lambda}(\mathbf{w}), \quad (1.100)$$

where the matrix $\tilde{\mathbb{T}}(\mathbf{w})$ has the vectors $\mathbf{r}_s(\mathbf{w})$, $s = 1, \dots, 4$ as its columns and

$$\tilde{\Lambda}(\mathbf{w}) = \text{diag} \left(\tilde{\lambda}_1(\mathbf{w}), \dots, \tilde{\lambda}_4(\mathbf{w}) \right). \quad (1.101)$$

This means that the matrix $\mathbb{A}_1(\mathbf{w})$ is diagonalizable.

Now we can show that the matrix $\mathbb{P}(\mathbf{w}, \mathbf{n})$ is diagonalizable too. By (1.90) and (1.100)

$$\mathbb{P}(\mathbf{w}, \mathbf{n}) = |\mathbf{n}| \mathbb{Q}^{-1}(\mathbf{n}) \tilde{\mathbb{T}}(\mathbb{Q}(\mathbf{n}) \mathbf{w}) \tilde{\Lambda}(\mathbb{Q}(\mathbf{n}) \mathbf{w}) \tilde{\mathbb{T}}(\mathbb{Q}(\mathbf{n}) \mathbf{w})^{-1} \mathbb{Q}(\mathbf{n}). \quad (1.102)$$

Under the notation

$$\begin{aligned} \Lambda(\mathbf{w}, \mathbf{n}) &= |\mathbf{n}| \tilde{\Lambda}(\mathbb{Q}(\mathbf{n}) \mathbf{w}) = \text{diag} |\mathbf{n}| \left(\tilde{\lambda}_1(\mathbb{Q}(\mathbf{n}) \mathbf{w}), \dots, \tilde{\lambda}_4(\mathbb{Q}(\mathbf{n}) \mathbf{w}) \right), \\ \mathbb{T}(\mathbf{w}, \mathbf{n}) &= \mathbb{Q}^{-1}(\mathbf{n}) \tilde{\mathbb{T}}(\mathbb{Q}(\mathbf{n}) \mathbf{w}), \end{aligned} \quad (1.103)$$

we see that the matrix $\mathbb{P}(\mathbf{w}, \mathbf{n})$ satisfies the relation $\mathbb{T}^{-1} \mathbb{P} \mathbb{T} = \Lambda(\mathbf{w}, \mathbf{n}) = \text{diag}(\lambda_1, \dots, \lambda_4)$ and thus \mathbb{T} diagonalizes the matrix \mathbb{P} .

Chapter 2

Computational methods in fluid dynamics

In this section we shall consider an initial-boundary value problem for the system of equations in the space-time cylinder $Q_T = \Omega \times (0, T)$, where $\Omega \in \mathbb{R}^2$. We give here a short overview of methods for discretization in space.

2.1 Triangulation

First step in the process of the discretization is to divide the computational domain $\overline{\Omega}$ into a finite number of subsets with properties described below. These subsets form the set, further denoted by \mathcal{T}_h , called the *triangulation of the domain* Ω . The parameter $h > 0$ of the triangulation usually represents maximum of diameters of all elements $K \in \mathcal{T}_h$. The elements $K \in \mathcal{T}_h$ are in the context of the finite volume method called *finite volumes*.

Properties of \mathcal{T}_h :

- a) Each $K \in \mathcal{T}_h$ is closed and connected with its interior $K^\circ \neq \emptyset$.
- b) Each $K \in \mathcal{T}_h$ has a Lipschitz boundary.
- c) $\cup_{K \in \mathcal{T}_h} K = \overline{\Omega}$
- d) If $K_1, K_2 \in \mathcal{T}_h$, $K_1 \neq K_2$, then $K_1^\circ \cap K_2^\circ = \emptyset$.

In our case of the two-dimensional problem, we assume that the domain Ω is obtained as an approximation of the original computational domain (also denoted by Ω), and the triangulation is chosen accordingly to the following attributes:

- A) Each $K \in \mathcal{T}_h$ is a closed triangle or quadrilateral, possibly with curved edges.

B) For $K_1, K_2 \in \mathcal{T}_h$, $K_1 \neq K_2$ we have either $K_1 \cap K_2 = \emptyset$ or K_1, K_2 share one edge (if the shared edge is a whole common edge, we call the triangulation *regular*) or K_1, K_2 share one vertex.

C) $\cup_{K \in \mathcal{T}_h} K = \overline{\Omega}$.

Furthermore

$$\mathcal{T}_h = \{K_i, i \in I\},$$

where $I \subset \mathbb{Z}^+ = \{0, 1, 2, \dots\}$ is a suitable index set.

By Γ_{ij} we denote a common edge between two neighboring elements K_i and K_j . We set

$$s(i) = \{j \in I; K_j \text{ is a neighbor of } K_i\}.$$

The boundary $\partial\Omega$ is formed by a finite number of faces of elements K_i adjacent to $\partial\Omega$. We denote all these boundary edges by S_j , where $j \in I_b \subset \mathbb{Z}^- = \{-1, -2, \dots\}$. Now we set

$$\gamma(i) = \{j \in I_b; S_j \text{ is a face of } K_i \in \mathcal{T}_h\}$$

and

$$\Gamma_{ij} = S_j \text{ for } K_i \in \mathcal{T}_h \text{ such that } S_j \subset \partial K_i, j \in I_b.$$

For K_i not containing any boundary face S_j we set $\gamma(i) = \emptyset$.

Obviously, $s(i) \cup \gamma(i) = \emptyset$ for all $i \in I$. If we write $S(i) = s(i) \cup \gamma(i)$, we have

$$\partial K_i = \cup_{j \in S(i)} \Gamma_{ij}, \quad \partial K_i \cap \partial\Omega = \cup_{j \in \gamma(i)} \Gamma_{ij}.$$

Model example Our model example in this chapter will be the scalar *linear advection-diffusion equation* with diffusion term left out:

$$\nabla \cdot \mathbf{f}(u(x_1, x_2)) = 0. \quad (2.1)$$

In particular, we will be concerned with the case

$$\mathbf{f}(u) = \beta u,$$

where $\beta(x_1, x_2) = (-x_2, x_1)/|\mathbf{x}|$ represents a circular counterclockwise flow field and $\Omega = [0, 1] \times [0, 1]$. We prescribe the following boundary condition:

$$u = g \text{ on } \Gamma_-, \text{ where} \quad (2.2)$$

$$g = 1 \text{ on } \Gamma_-^1, \text{ and} \quad (2.3)$$

$$g = 0 \text{ on } \Gamma_-^2, \quad (2.4)$$

where $\Gamma_- = \{\mathbf{x} = (x_1, x_2) : \beta(x_1, x_2) \cdot \mathbf{n}(\mathbf{x}) < 0\}$, $\Gamma_1 = [0, 0.5] \times 0$, and $\Gamma_-^2 = \Gamma_- \setminus \Gamma_-^1$. We do not prescribe any boundary condition on the rest of the boundary, i.e. on $\Gamma = \partial\Omega/\Gamma_-$.

2.2 Finite Volume method

2.2.1 Scheme derivation

Let us assume that $u : \bar{\Omega} \rightarrow \mathbb{R}$ is a classical (\mathcal{C}^1) solution of the equation (2.1), $\mathcal{T}_h = \{K_i\}_{i \in I}$ is the finite volume mesh as an approximation of Ω . Integrating the equation (2.1) over K_i and using Green's theorem on K_i , we get the identity

$$\int_{\partial K_i} \mathbf{f}(u) \cdot \mathbf{n} dS = 0. \quad (2.5)$$

Now we shall approximate the integral averages $\int_{K_i} u(\mathbf{x}) d\mathbf{x} / |K_i|$ of the quantity u over the finite volume K_i by u_i :

$$u_i \approx \frac{1}{|K_i|} \int_{K_i} u(\mathbf{x}) d\mathbf{x}, \quad (2.6)$$

called the value of *the approximate solution* on K_i . Further, we approximate the flux $\mathbf{f}(u)(\mathbf{n}_{ij})$ of the quantity u through the edge Γ_{ij} in the direction \mathbf{n}_{ij} with the aid of a *numerical flux* $\mathbf{H}(u_i, u_j, \mathbf{n}_{ij})$, depending on the value of the approximate solution u_i on the finite volume K_i , the value u_j on K_j , and on the normal \mathbf{n}_{ij} . If $\Gamma_{ij} \subset \partial\Omega$ (i.e. the finite volume K_i is adjacent to $\partial\Omega$, $j \in \gamma(i)$), then there is no neighbor K_j of K_i adjacent to the face Γ_{ij} from the exterior of Γ , and it is necessary to specify u_j on the basis of boundary conditions. This will be explained further in the section about the discretization of the Euler equations.

The scheme is also equipped with *initial conditions* u_i^0 , $i \in I$, defined by

$$u_i^0 = \frac{1}{|K_i|} \int_{K_i} u^0(\mathbf{x}) d\mathbf{x}, \quad (2.7)$$

under the assumption that the function u^0 from (??) is locally integrable, i.e. that $u^0 \in L^1_{\text{loc}}(\Omega)$.

2.2.2 Properties of the numerical flux

In what follows, we shall assume that the numerical flux \mathbf{H} has the following properties:

- A) $\mathbf{H}(u, v, \mathbf{n})$ is defined and continuous on $\mathcal{D} \times \mathcal{D} \times \mathcal{S}_1$, where \mathcal{D} is the domain of definition of the flux \mathbf{f} and \mathcal{S}_1 is the unit sphere in \mathbb{R}^N (here $N = 2$).
- B) \mathbf{H} is *consistent*:

$$\mathbf{H}(u, u, \mathbf{n}) = \mathbf{f}(u) \mathbf{n}, \quad u \in \mathcal{D}, \quad \mathbf{n} \in \mathcal{S}_1.$$

C) \mathbf{H} is conservative:

$$\mathbf{H}(u, v, \mathbf{n}) = -\mathbf{H}(v, u, -\mathbf{n}), \quad u, v \in \mathcal{D}, \quad \mathbf{n} \in \mathcal{S}_1.$$

We define a finite volume approximate solution of (2.1) as piecewise constant function u_h , defined a.e. in Ω so that $u_h|_{K_i^o} = u_i$ for all $i \in I$, where K_i^o is the interior of K_i , i.e. $K_i^o = K_i \setminus \partial K_i$, and u_i are obtained from the formula

$$\frac{1}{K_i} \sum_{j \in S(i)} \mathbf{H}(u_i, u_j, \mathbf{n}_{ij}) = 0, \quad K_i \in \mathcal{T}_h. \quad (2.8)$$

The number u_i is the value of the approximate solution on the finite volume K_i .

2.2.3 Example

Now we give an example of a numerical flux, which is the *upwinding* numerical flux defined by

$$H(u_i, u_j, \mathbf{n}_{ij}) = \begin{cases} u_i, & \text{if } \beta \cdot \mathbf{n}_{ij} > 0 \\ u_j, & \text{if } \beta \cdot \mathbf{n}_{ij} \leq 0 \end{cases}. \quad (2.9)$$

Solutions of the model example using the finite volume method with the numerical flux (2.9) are presented below. The computation was obtained on three different meshes.

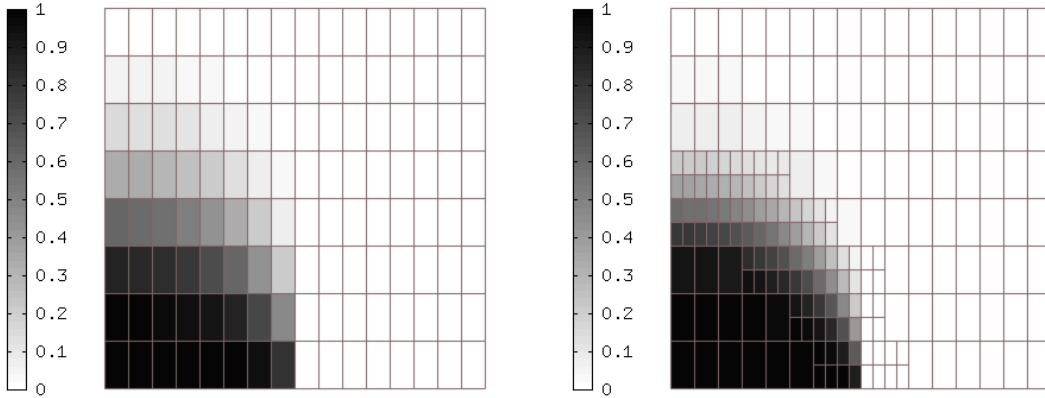


Figure 2.1: Solution u on a mesh with 128 (left), and 206(right) elements.

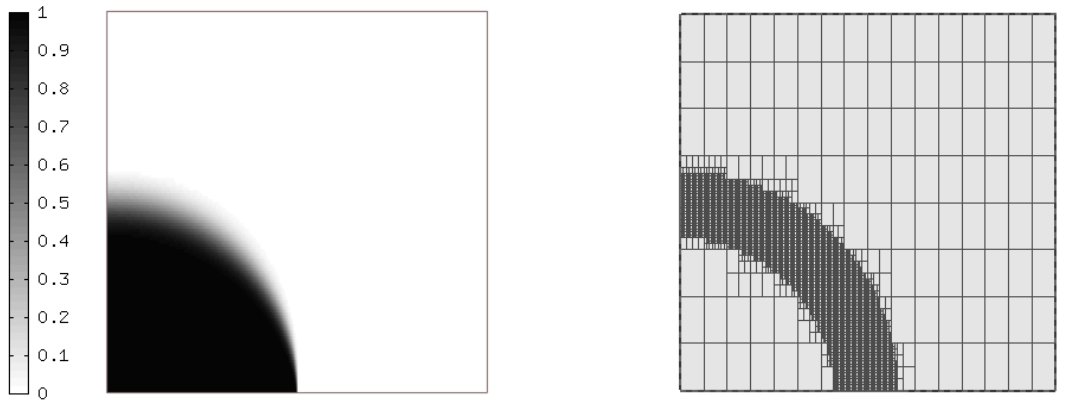


Figure 2.2: Solution u (left) and the corresponding mesh (right) with 4751 elements.

We can see that the solution is smeared near the steep gradient. We would like to avoid this.

2.3 Finite Element Method

Application of the standard FEM to the problem (2.1) is very straightforward. We multiply the equation (2.1) with a test function v from a suitable finite-dimensional approximation of the $H^1(\Omega)$ space. We then integrate over Ω and use the Green's theorem. We get the following identity:

$$\int_{\Omega} \mathbf{f}(u) \cdot \nabla v d\mathbf{x} = \int_{\Gamma_-} g v dS, \quad (2.10)$$

which must hold for all v from the finite-dimensional approximation of the $H^1(\Omega)$ space. As this finite-dimensional sub-space of the $H^1(\Omega)$ space, we usually take the space of continuous functions, which are polynomials up to order p on each element of the triangulation.

For singularly perturbed equations, or strictly first order equations, such as our model equation, the standard Galerkin FEM exhibits *Gibbs phenomenon*, manifested by *spurious oscillations* in the numerical solution. This behavior is well noticeable in the following figures.

2.3.1 Example

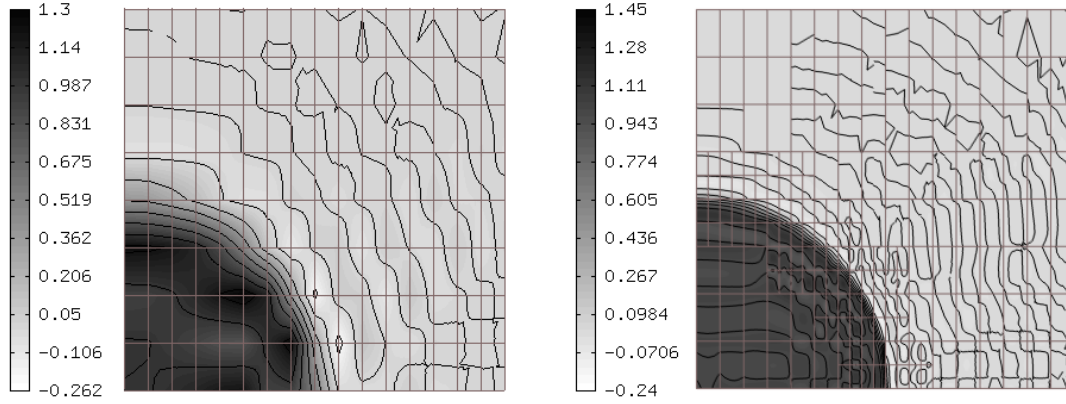


Figure 2.3: Solution u on a mesh with 128 elements and $p = 1$ (left), and 206 elements and $p = 2$ (right).

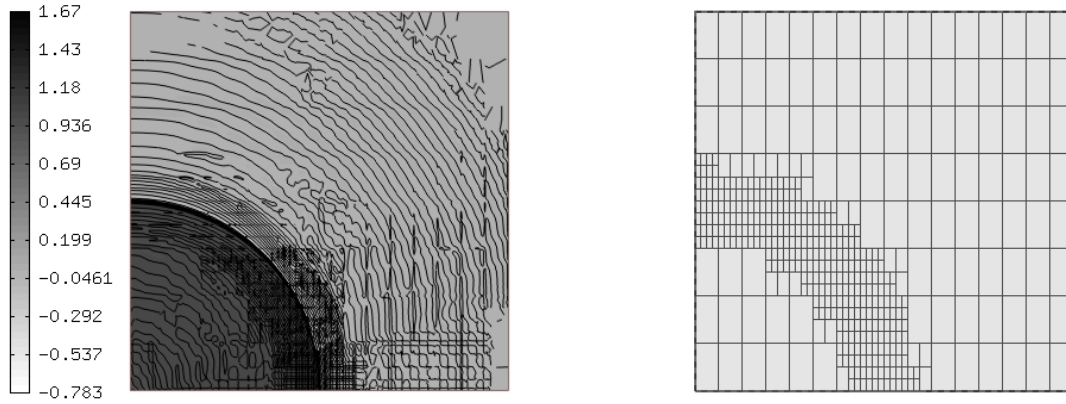


Figure 2.4: Solution u (left) and the corresponding mesh (right) with 449 elements, $p = 5$, and 11206 degrees of freedom.

The oscillations can be diminished using stabilization technique, e.g. streamwind-upwind petrov-galerkin stabilization. Their disadvantage is the necessity to fine tune various parameters for them to function properly. For a reference, see e.g. [?], [?], [?].

2.4 Discontinuous Galerkin method

For complex problems of compressible flow, the choice of optimal stabilization parameters becomes quite sophisticated and difficult. Due to this reason, there was an effort to develop methods which would not need such stabilization techniques, and

would still offer reasonable resolution of shockwaves, boundary and interior layers, and steep gradients without exhibiting spurious oscillations in an approximate solution. It is based on the idea to combine finite volume and finite element methods leading to the so-called *discontinuous Galerkin finite element method (DGFEM)*. Here we shall derive and analyze DGFEM for our model example. Let \mathcal{T}_h be a triangulation of Ω . For each $K \in \mathcal{T}_h$ we introduce the notation

$$\partial K^- = \{x \in \partial K; \beta(\mathbf{x}) \cdot \mathbf{n}(\mathbf{x}) < 0\}, \quad (2.11)$$

$$\partial K^+ = \{x \in \partial K; \beta(\mathbf{x}) \cdot \mathbf{n}(\mathbf{x}) \geq 0\}. \quad (2.12)$$

By $H^k(\Omega, \mathcal{T}_h)$ we denote the so-called *broken Sobolev space*:

$$H^k(\Omega, \mathcal{T}_h) = \{v \in L^2(\Omega); v|_K \in H^k(K) \forall K \in \mathcal{T}_h\}. \quad (2.13)$$

For $u \in H^1(\Omega, \mathcal{T}_h)$ we set

$$u_K^+ = \text{trace of } u|_K \text{ on } \partial K \quad (2.14)$$

(i.e. the interior trace of u on ∂K). For each edge $E \subset \partial K \setminus \Gamma$ of K , there exists $K' \neq K$, $K' \in \mathcal{T}_h$, adjacent to E from the opposite side than K . Then we put

$$u_K^- = \text{trace of } u|_{K'} \text{ on } E. \quad (2.15)$$

In this way we obtain the exterior trace u_K^- of u on $\partial K \setminus \Gamma$ and define the jump of u on $\partial K \setminus \Gamma$:

$$[u]_K = u_K^+ - u_K^-. \quad (2.16)$$

2.4.1 Derivation of the DGFEM

Let $u \in H^1(\Omega)$ be a solution of the problem (2.1). Then u satisfies the identity

$$\int_K \nabla \cdot (\beta u) \varphi d\mathbf{x} = 0, \quad \varphi \in H^1(\Omega, \mathcal{T}_h), \quad K \in \mathcal{T}_h. \quad (2.17)$$

The application of Green's theorem gives

$$\begin{aligned} \int_K \nabla \cdot (\beta u) \varphi d\mathbf{x} &= \int_{\partial K} (\beta u_K^+) \cdot \mathbf{n} \varphi_K^+ dS - \int_K (\beta u) \cdot \nabla \varphi d\mathbf{x} \\ &= \int_{\partial K^-} (\beta u_K^+) \cdot \mathbf{n} \varphi_K^+ dS + \int_{\partial K^+} (\beta u_K^+) \cdot \mathbf{n} \varphi_K^+ dS \\ &\quad - \int_K (\beta u) \cdot \nabla \varphi d\mathbf{x}, \end{aligned} \quad (2.18)$$

where \mathbf{n} is the unit outer normal to the element boundary ∂K . As $u \in H^1(\Omega)$, we have $u_K^- = u_K^+$. Moreover, $u_K^-|_{\partial K \cap \Gamma} := u|_{\partial K \cap \Gamma} = g$. Then we can write

$$\begin{aligned} \int_K \nabla \cdot (\beta u) \varphi d\mathbf{x} &= \int_{\partial K^-} (\beta u_K^-) \cdot \mathbf{n} \varphi_K^+ dS + \int_{\partial K^+} (\beta u_K^+) \cdot \mathbf{n} \varphi_K^+ dS \\ &\quad - \int_K (\beta u) \cdot \nabla \varphi d\mathbf{x}. \end{aligned} \quad (2.19)$$

Applying Green's theorem again, we get the identity

$$\int_K \nabla \cdot (\beta u) \varphi d\mathbf{x} = \int_{\partial K^-} (\beta (u_K^+ - u_K^-)) \cdot \mathbf{n} \varphi_K^+ dS. \quad (2.20)$$

Setting

$$\begin{aligned} a_K(u, \varphi) &= \int_K \nabla \cdot (\beta u) \varphi d\mathbf{x} - \int_{\partial K^- \setminus \Gamma} (\beta [u]_K) \cdot \mathbf{n} \varphi_K^+ dS \\ &\quad - \int_{\partial K^- \cap \Gamma_-} (\beta u_K^+) \cdot \mathbf{n} \varphi_K^+ dS, \end{aligned} \quad (2.21)$$

$$L_K(\varphi) = \int_{\partial K^- \cap \Gamma_-} (\beta g) \cdot \mathbf{n} \varphi_K^+ dS, \quad (2.22)$$

we can rewrite the equation (2.20) as

$$a_K(u, \varphi) = L_K(\varphi), \quad \varphi \in H^1(K), \quad K \in \mathcal{T}_h. \quad (2.23)$$

This identity makes sense also for $u \in H^1(\Omega, \mathcal{T}_h)$. In this case, we can note that in (2.18), on ∂K^- (= the inlet of K with respect to the velocity β) we replace the value u_K^+ (the interior trace of u) by u_K^- . This means that *upwinding* is used here, because the value of the trace of u on ∂K^- is taken from the side of ∂K^- against the velocity direction.

Now, on the basis of (2.23) we come to the definition of the *discrete problem*. The *approximate solution* is a function u_h satisfying the conditions

$$\text{A) } u_h \in S_h = S^{p,-1}(\Omega, \mathcal{T}_h) := \{\varphi \in L^2(\Omega); \varphi|_K \in P^p(K) \forall K \in \mathcal{T}_h\}$$

$$\text{B) } a_K(u_h, \varphi_h) = L_K(\varphi_h), \quad \forall \varphi_h \in S_h, \quad \forall K \in \mathcal{T}_h.$$

Here, p is an integer. In general on each element a different polynomial degree can be used for the approximation. The approximate solution and test functions are piecewise polynomial functions without any continuity requirement on interfaces between neighboring elements. The continuity requirement is replaced here by the jump term $\int_{\partial K^- \setminus \Gamma} (\beta [u]_K) \cdot \mathbf{n} \varphi_K^+ dS$.

We can see that we successfully avoided the smeared layer close to the steep gradient, and also the Gibbs phenomenon in the whole domain. What we have to deal here with on the other hand are overshoots and undershoots close to the step gradient. There is several approaches, some of which are described in the last chapter with numerical experiments, for the case of the Euler equations.

2.4.2 Example

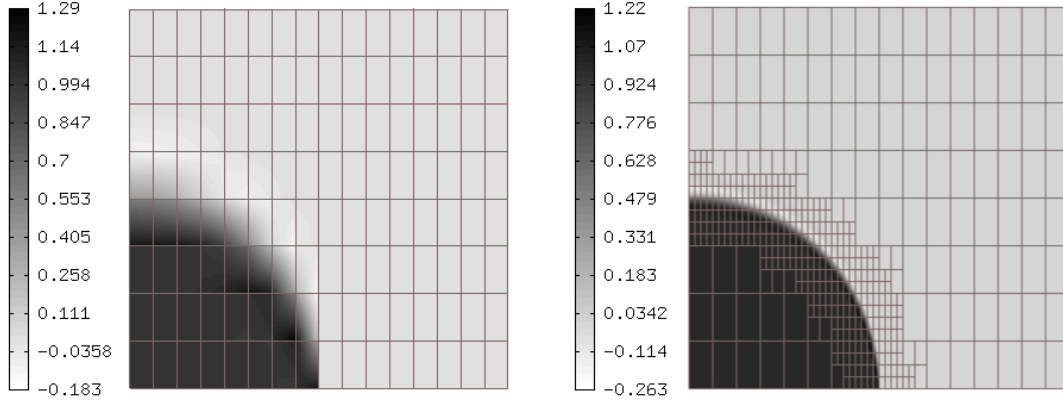


Figure 2.5: Solution u on a mesh with 128 elements and $p = 1$ (left), and 206 elements and $p = 2$ (right).

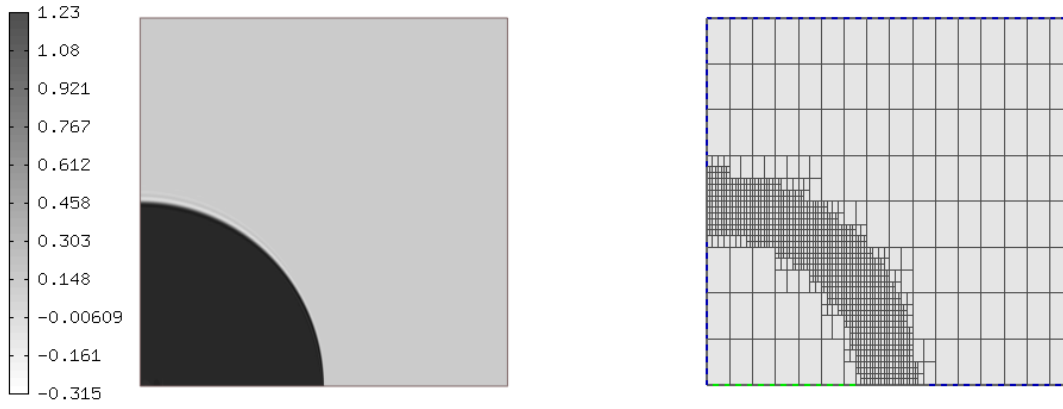


Figure 2.6: Solution u (left) and the corresponding mesh (right) with 1271 elements, $p = 3$, and 20336 degrees of freedom.

Chapter 3

Adaptive algorithms

In computational fluid dynamics we are interested in the computation of sufficiently accurate solutions of various flow problems. The goal of the design of any numerical method is *reliability* and *efficiency*. Reliability means that the computational error is controlled at a given tolerance level. The efficiency means that the cost of the computation of a solution within a given tolerance is as small as possible. These two requirements are usually achieved with the aid of mesh refinement techniques. The goal is to achieve reliability either in the sense that the numerical solution approximates the exact solution in a given norm within a given tolerance, or in the sense that some physically relevant quantities (e.g. flux through a part of a boundary, drag, lift) are computed within a given tolerance.

Standard h-adaptive FEM, where (adaptive) mesh refinement is used, is a well known, and well spread method. Sometimes, much faster convergence can be achieved by increasing the polynomial degree of the elements instead (*p*-refinement). Such approach is far more efficient for elements where the solution is smooth. In the following section, a brief description of this technique is presented.

3.1 p -adaptivity

In the process of FEM, or DGFEM usage, the domain Ω is first approximated by a suitable (often polygonal) domain Ω_h . This is one of the so-called *variational crimes* – departures from the “mathematically clean” variational framework – since $\Omega_h \not\subset \Omega$ and the solution or other functions from the weak formulation are not defined where they are to be approximated or evaluated. In practice these crimes are often hard to avoid. In this work, the domain Ω , and its approximation used in computations Ω_h will be both denoted Ω .

3.1.1 Higher-order finite element space

Let the domain Ω be covered with a mesh $\mathcal{T}_h = \{K_1, K_2, \dots, K_M\}$ where the elements K_m carry arbitrary polynomial degrees $1 \leq p_m, m = 1, 2, \dots, M$. The broken Sobolev space $H^1(\Omega, \mathcal{T}_h)$ is now approximated by a space of piecewise-polynomial functions

$$V_h = \{v \in L^2(\Omega); v|_{K_m} \in P^{p_m}(K_m) \text{ for all } 1 \leq m \leq M\}$$

where P^p is defined as

$$P^p = \text{span}\left\{ \sum_{0 \leq i, j \leq p} \alpha_i x_1^i x_2^j, \alpha_i \in \mathbb{R} \right\}$$

in case of triangles and

$$P^p = \text{span}\left\{ \sum_{\substack{0 \leq i, j \leq p \\ i+j \leq p}} \alpha_i x_1^i x_2^j, \alpha_i \in \mathbb{R} \right\}$$

in case of quadrilaterals.

In discontinuous Galerkin method, there is no difference between basis functions that logically belong to vertices, edges, or volumes. They are all treated as basis functions that belong to a vertex volume. Also the name (vertex-, edge- basis functions) is dropped and only the generic term *basis functions* is used. The next section shows one possible way of constructing basis functions for the discontinuous Galerkin method.

3.1.2 Hierarchic shape functions

Hierarchic shape functions are constructed in such a way that the basis \mathcal{B}^{p+1} of the polynomial space $P^{p+1}(K)$ is obtained from the basis \mathcal{B}^p of the polynomial space $P^p(K)$ by adding new shape functions only. This is essential for p - and hp -adaptive finite element codes since one does not have to change the shape functions

completely when increasing the degree of polynomial approximation. In this section we will describe the popular *Lagrange-based* hierarchic shape functions.

Definition 3.1 (Legendre polynomials) *We define the Legendre polynomials of degree k as*

$$L_k(x) = \frac{1}{2^k k!} \frac{d^k}{dx^k} (x^2 - 1)^k, \quad k = 0, 1, 2, \dots$$

The set of Legendre polynomials forms an orthonormal basis of the space $L^2(-1, 1)$.

Quadrilaterals On quadrilaterals, the set of shape functions is the product of Legendre polynomials. Displayed are basis functions up to the degree 3.

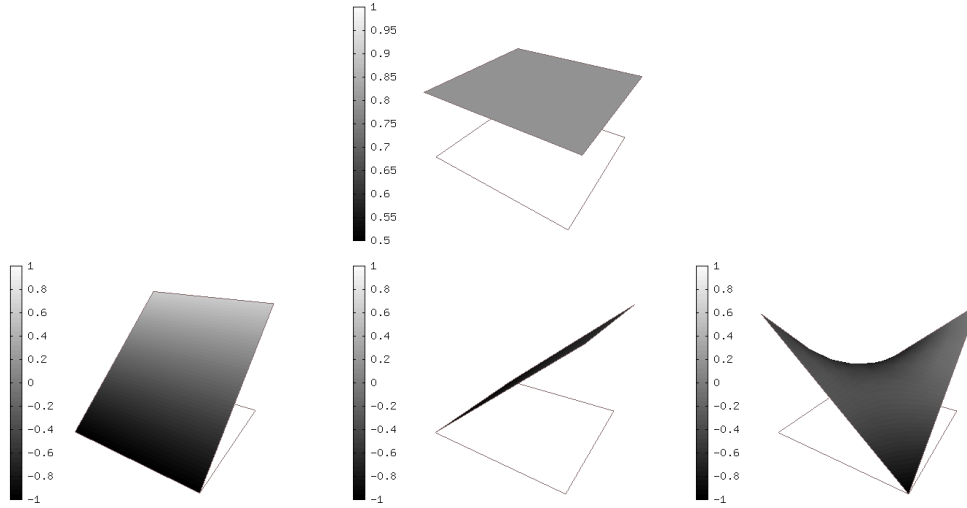


Figure 3.1: Shape functions - degree 0 top, degree 1 bottom

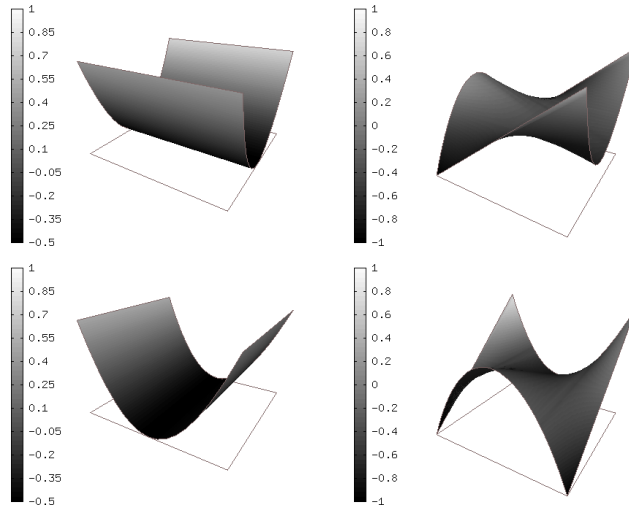


Figure 3.2: Shape functions - degree 2

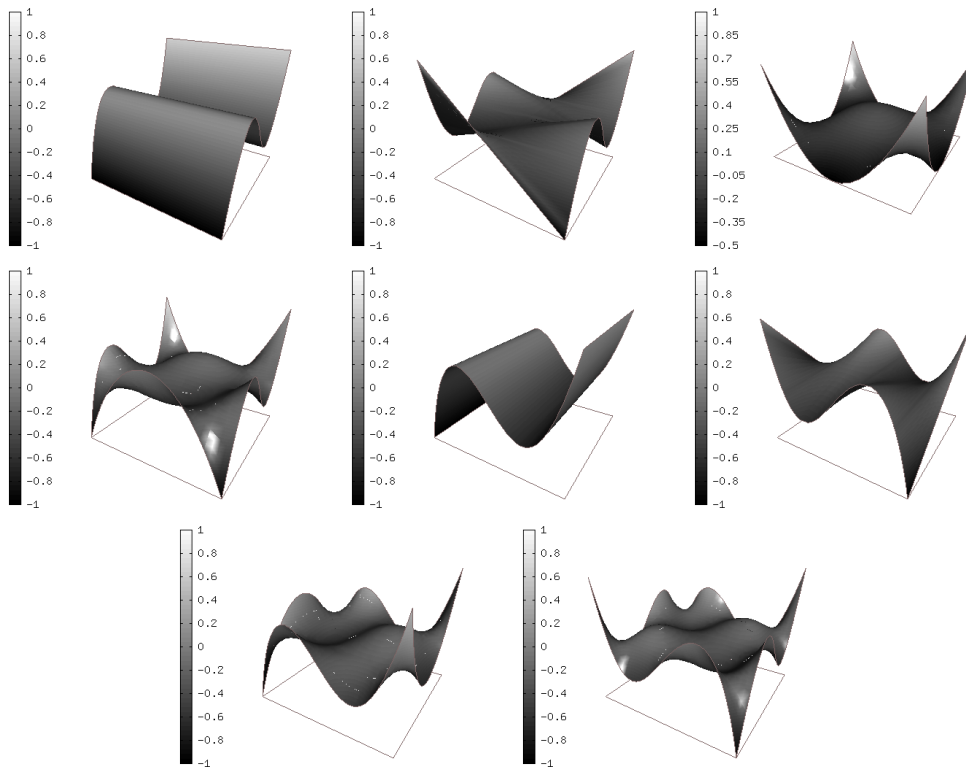


Figure 3.3: Shape functions - degree 3

Triangles Situation on triangles is a bit more complicated. Since the number of basis functions spanning the space of polynomials of certain degree is not the same

as on quadrilaterals, the set of shape functions can not be the product of Legendre polynomials. In the construction of the set, the following definition will be used.

Definition 3.2 (Affine coordinates) Barycentric coordinates *on the triangular reference domain \hat{K}* (Figure ??) are defined as

$$\lambda_1(\xi_1, \xi_2) = \frac{\xi_2 + 1}{2}, \quad \lambda_2(\xi_1, \xi_2) = -\frac{\xi_1 + \xi_2}{2}, \quad \lambda_3(\xi_1, \xi_2) = \frac{\xi_1 + 1}{2}$$

Now the basis functions are defined in the following way:

$$L_{ij}^{tri} = L_i(\lambda_3(x, y) - \lambda_2(x, y)) \cdot L_j(\lambda_2(x, y) - \lambda_1(x, y)), \quad (3.1)$$

the resulting set is then $L_p^{tri} = \{L_{ij}^{tri}, i, j = 0, \dots, p\}$, where p represents the highest polynomial degree we use for approximation. Figures of the set for $p = 3$ follows.

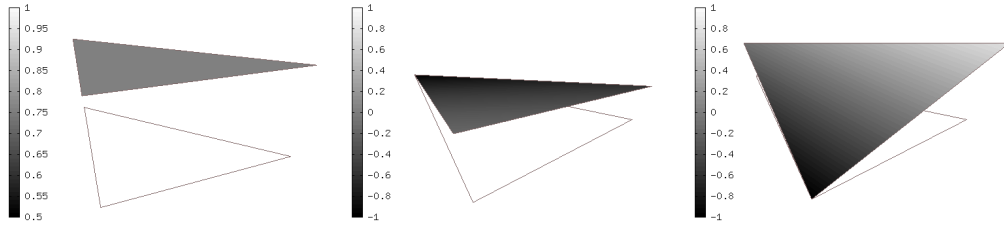


Figure 3.4: Shape functions - degree 0 on the left, degree 1 rest

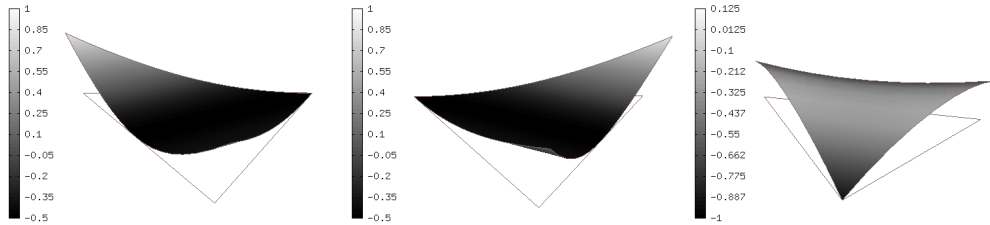


Figure 3.5: Shape functions - degree 2

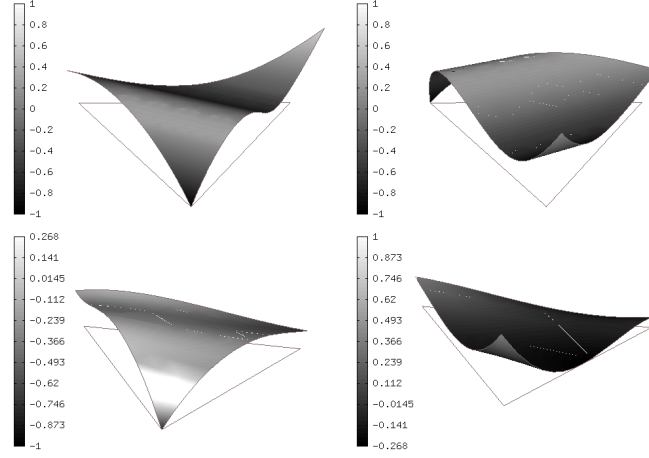


Figure 3.6: Shape functions - degree 3

Proposition 3.1 *The shape functions described above, both on quadrilaterals, and on triangles constitute a hierarchic basis of the space $P^m(\hat{K})$*

Proof The complete proof can be found in [?]. Briefly, the following must be verified:

- A) Are all the shape functions linearly independent?
- B) Do they all belong to the space $P^m(\hat{K})$?
- C) Matches their number the dimension of the space $P^m(\hat{K})$?

3.1.3 Higher-order numerical quadrature

Most commonly, the integrals in (??), (??) are evaluated numerically by the *Gaussian quadrature*. The k -point Gaussian quadrature rule on the domain \hat{K} has the form

$$\int_{\hat{K}} g(\xi) d\xi \approx \sum_{i=1}^k w_{k,i} g(\xi_{k,i}) \quad (3.2)$$

where g is a bounded continuous function, $\xi_{k,i} \in \hat{K}, i = 1, 2, \dots, k$ are the integration points and $w_{k,i} \in \mathbb{R}$ are the integration weights. The sum of the integration weights must be equal to the area of \hat{K} , so that the rule (3.2) is exact for constants. If the points and weights are chosen carefully, the formula (3.2) can be exact for polynomials up to a certain degree $q > 0$.

In 1D the integration points are roots of the Legendre polynomials. Also in 2D it is quite easy to find the integration points and weights for low-degree polynomials occurring in traditional linear FEM. For higher-degree polynomials, however, the task of finding optimal Gaussian quadrature rules presents a complex non-linear

optimization problem with many open questions left. Optimal integration points and weights are known on \hat{K} for polynomials up to degree 10. Suboptimal (with more points than necessary) rules have been found for polynomials up to degree 20. Complete integration tables along with more information on this subject can be found in [?].

3.2 Hanging nodes

As the goal is automatic mesh adaptivity on both triangular and quadrilateral (or combined) meshes, one of the important characteristics of triangulations is their regularity. This section shows why the regularity assumption should be dropped and how irregular meshes are treated.

3.2.1 Hanging nodes and irregularity rules

At the beginning, let us recall the *red-green* refinement strategy. This technique first subdivides desired elements into geometrically convenient subelements with hanging nodes and then it eliminates the hanging nodes by forcing refinement of additional elements. This is illustrated in Figure 3.7. This approach preserves regularity of the mesh but in the case of triangles, it introduces elements with sharp

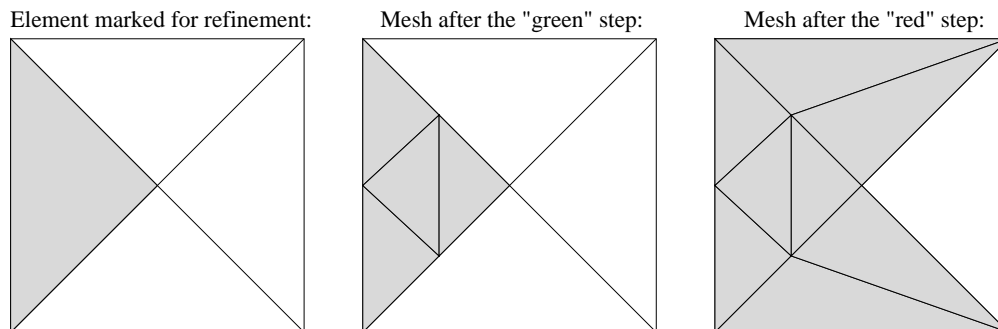


Figure 3.7: Red-green refinement.

angles which are not desirable in finite element analysis. In the case of quadrilaterals, it also negatively influence the the element shapes. It becomes extremely cumbersome when repeated refinements occur in the same part of the domain (e.g., toward a boundary layer or point singularity).

The “red” refinements can be avoided by introducing *hanging nodes*, i.e., by allowing *irregular meshes* where element vertices lie in the interior of edges of other elements. In order to keep the computer implementation simple, most finite element working with hanging nodes limit the maximum difference of refinement levels of adjacent elements to one (so-called *1-irregularity rule*) – see, e.g., [?, ?, ?]. In the following, by *k-irregularity rule* (or *hanging nodes of level k*) we mean this type of restriction where the maximum difference of refinement levels of adjacent elements is k . In this context, $k = 0$ corresponds to adaptivity with regular meshes and $k = \infty$ to adaptivity with arbitrary-level hanging nodes.

It is illustrated in Figure 3.8 that even the 1-irregularity rule does not avoid all forced refinements:

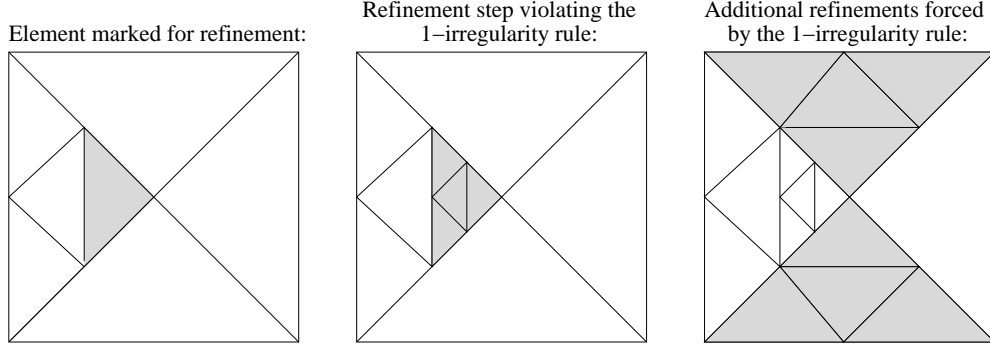


Figure 3.8: Refinement with 1-irregularity rule.

The amount of forced refinements in the mesh depends strongly on the level of hanging nodes. Next we introduce a model problem and show that the level of hanging nodes also influences significantly both the number of degrees of freedom and condition number of stiffness matrices.

Model problem Consider a Poisson equation $-\Delta u = f$ in Ω with $u = 0$ on the boundary of Ω , where $\Omega = (-1, 1)^2$. Assume a right-hand side f such that the corresponding exact solution u is zero everywhere in Ω with the exception of a local perturbation occurring inside of a small triangle T_n with the vertices $[-2^{-n}, -2^{-n}]$, $[0, 0]$, $[-2^{-n}, 2^{-n}]$. The domain is covered with a coarse four-element mesh shown in Figure 3.9. For simplicity, all elements are equipped with a uniform polynomial degree $p \geq 1$.

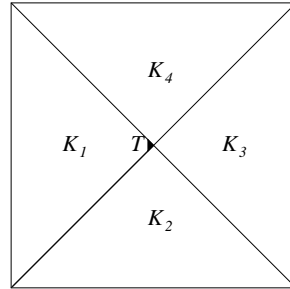


Figure 3.9: Domain and initial mesh for the study of k -irregularity rules.

We perform the following simple experiment: Starting from the coarse mesh shown in Figure 3.9, we run an adaptive algorithm which in each step applies the “green” refinement to every mesh triangle K such that $K \cap T_n \neq \emptyset$ and $K \not\subset T_n$. Figure 3.10 shows final meshes obtained with various levels of hanging nodes for $n = 5$:

Notice that, since $T_n \subset K_1$, all refinements within the elements K_2, K_3, K_4 are forced and they do not improve the resolution. The reader can see that the amount

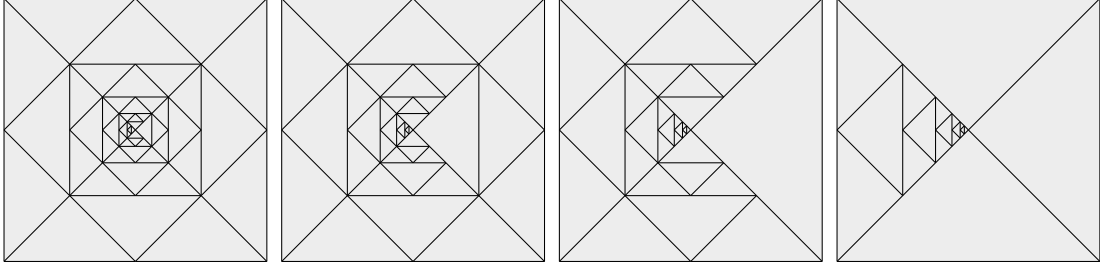


Figure 3.10: Meshes obtained with k -irregularity rules, $k = 1, 2, 3, \infty$.

of forced refinements decreases as k grows, vanishing completely with $k = \infty$.

Next let us choose, e.g., $p = 3$, and run the same adaptive procedure for $n = 1, 2, \dots, 15$. Figure 3.11 shows the number of degrees of freedom corresponding to the final meshes. The horizontal axis represents the spatial scale 2^{-n} . For the same computations, Figure 3.12 shows the condition number of the corresponding stiffness matrices.

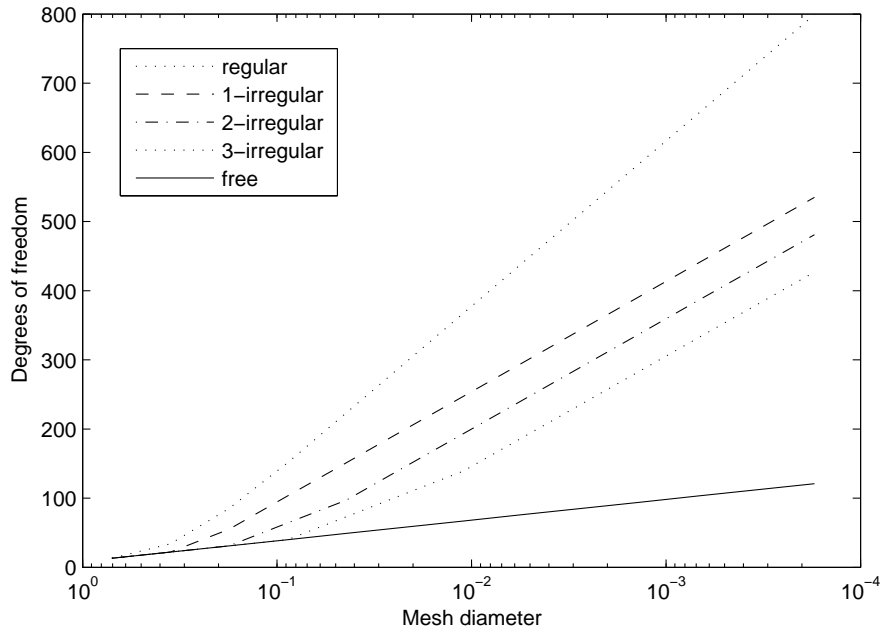


Figure 3.11: Size of the stiffness matrix for final meshes. Level of hanging nodes: $k = 0, 1, 2, 3, \infty$.

In DGFEM, the requirement from the conforming FEM of continuity across mesh edges is dropped, and this means that we only have to be able to calculate fluxes \mathbf{f} across mesh edges properly (as opposed to conforming finite elements, where we have to take care of the continuity requirement in vertices and on edges).

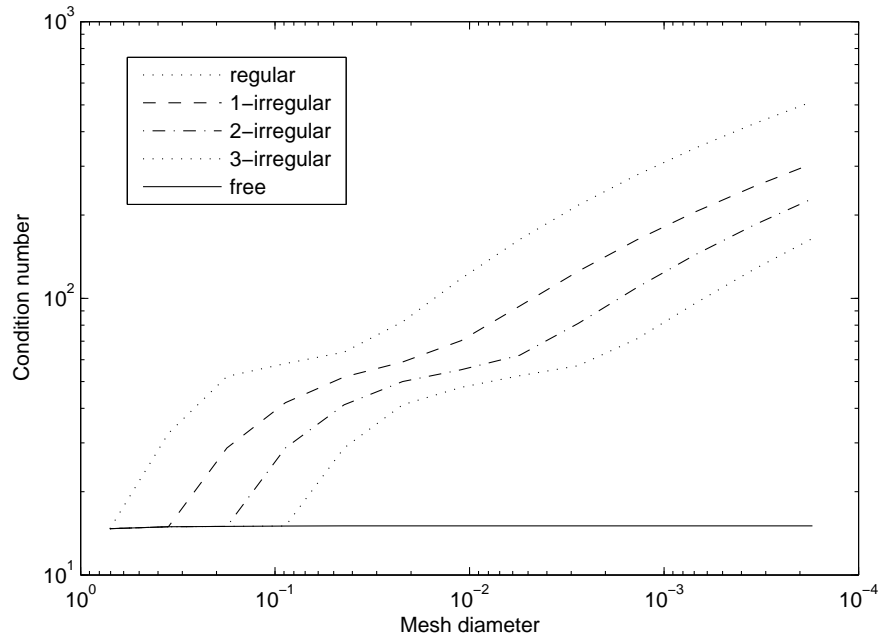


Figure 3.12: Condition number of the stiffness matrix for final meshes. Level of hanging nodes: $k = 0, 1, 2, 3, \infty$.

With arbitrary-level hanging nodes, the flux calculation also poses a challenge, and the algorithmic treatment will be described in the section 3.3.

3.3 Data structures for *hp*-adaptivity

The need for mesh adaptivity brings further demands that must be taken into account in the design of the data structures. A tree-like hierarchy of element refinement levels must be introduced and also more complicated neighbor search algorithms are required to correctly initialize the edge and vertex nodes and the integration points on the edges. The complexity of the data structures reflects strongly on the complexity of the rest of the finite element code.

Simplifying assumptions, such as the k -irregularity rules, are often imposed on the mesh with the aim of reducing code complexity. However, these assumptions can deteriorate the numerical performance of the code.

In the following we present an original design of data structures for adaptive meshes, free of any regularity assumptions yet retainign extreme simplicity.

Calculation of fluxes in DGFEM with arbitrary-level hanging will also be described.

3.3.1 Node and element structures

It should be noted first that we have departed from traditional mesh data structures storing complete information about the finite element problem. Our mesh structures contain geometrical information only. The remaining data, including basis function indices, boundary conditions, polynomial degrees, etc., are stored in separate structures describing concrete finite element spaces (H^1 , $\mathbf{H}(\text{curl})$, ...) and are accessible via the identification numbers of nodes and elements. This was done to allow multiple spaces and multiple element types to co-exist on the same mesh, which is vital for solving multi-physics and coupled problems. See [?]

The entire mesh is defined by two arrays of the following C structures. The first one stores all nodes:

```
structure Node
{
    // Type of the Node — either vertex , or edge.
    type;

    // Identifier of the particular Node.
    int id;

    // Physical coordinates of this Node, in case it is a vertex Node.
    coordinate x;
    coordinate y;

    // Two neighboring elements of this Node, in case it is an edge Node.
    Element elements[2];
};
```

This structure defines both vertex and edge nodes by utilizig the construct to distinguish between an edge and a vertex node. The standard vertex positions \mathbf{x} , \mathbf{y} , while typically stored separately, were placed in the vertex variant of the **Node** structure for simplicity. The edge variant contains pointers to the two elements

sharing the edge node (useful eg. when looking for a neighbor of an element sharing a particular edge in order to calculate the numerical flux across the edge).

The member **type** determines the variant of the structure. The remaining members are omitted for simplicity.

Elements are defined by the second structure:

```
struct Element
{
    // Type of the Element — either triangle, or quadrilateral.
    type;

    // Identifier of the particular Element.
    int id;

    // Info whether this element is active (has not been refined further).
    boolean active;

    // vertex nodes
    Node vn[4];

    // edge nodes
    Node en[4];
};
```

An element can be either active or inactive, hence the one-bit variable **active**. Active elements store pointers to the appropriate vertex and edge nodes. Inactive elements are those which have been split and are excluded from the computation. Their purpose is to hold pointers to the descendant elements.

Triangular and quadrilateral elements share the same structure and are distinguished by the member **type**. The rest of the variables are analogous to the **Node** structure.

3.3.2 Eliminating neighbor search by hashing

To properly initialize edge node pointers after reading a mesh file, one has to construct neighbor lists for all elements and use them in such a way that only one node is created for each edge. Further problems arise when certain elements are refined after automatic mesh adaptation. Unless hanging nodes are removed by extra refinements, no longer is each edge shared by at most two elements (there can be more of them). Standard neighbor lists fail to fully capture such situations and thus more complex data structures, eg. trees [?], have to be introduced.

Loosely inspired by [?], we have avoided all of these problems by introducing a function which, given the **id** numbers of two nodes, returns a pointer to the node halfway between them. If no such node exists, it is created first. The task of translating two numbers to a node pointer is accomplished using a hash table.

We are maintaining two independent layers of nodes: the first layer contains all vertex nodes, the second all edge nodes. The following two functions can be called:

```
Node* get_vertex_node(int id1, int id2);
Node* get_edge_node(int id1, int id2);
```

All nodes, apart from being accessible by their `id` number, can be reached using these functions by passing the `ids` of their “parent” nodes. Top-level vertex nodes (those loaded from the mesh file) are stored at the beginning of the node array and can be accessed directly without hashing. Mesh initialization then becomes trivial:

```

for all Elements
    vertex_numbers = read_element_vertex_id_numbers
    for all edges of the Element
        Element.vn = vertex_numbers
        Element.en[0] = get_edge_node(vertex_numbers[0], vertex_numbers[1]);
        Element.en[1] = get_edge_node(vertex_numbers[1], vertex_numbers[2]);
        Element.en[2] = get_edge_node(vertex_numbers[2], vertex_numbers[3]);
        // In case of triangles
        Element.en[3] = get_edge_node(vertex_numbers[3], vertex_numbers[0]);
    }
}

```

Element refinement is also very straightforward. No care must be taken of the neighboring elements, regardless of their refinement level or even existence:

```

create_triangle(Node v0, Node v1, Node v2)
{
    Element e;
    e.vn[0] = v0;
    e.vn[1] = v1;
    e.vn[2] = v2;
    e.en[0] = get_edge_node(v0.id, v1.id);
    e.en[1] = get_edge_node(v1.id, v2.id);
    e.en[2] = get_edge_node(v2.id, v0.id);
}

void refine_triangle(Element e)
{
    Node x0 = get_vertex_node(e.vn[0].id, e.vn[1].id);
    Node x1 = get_vertex_node(e.vn[1].id, e.vn[2].id);
    Node x2 = get_vertex_node(e.vn[2].id, e.vn[0].id);
    e.sons[0] = create_triangle(e.vn[0], x0, x2);
    e.sons[1] = create_triangle(x0, e.vn[1], x1);
    e.sons[2] = create_triangle(x2, x1, e.vn[2]);
    e.sons[3] = create_triangle(x0, x1, x2);
    e.active = 0;
}

create_quad(Node v0, Node v1, Node v2)
{
    Element e;
    e.vn[0] = v0;
    e.vn[1] = v1;
    e.vn[2] = v2;
    e.vn[3] = v3;

    e.en[0] = get_edge_node(v0.id, v1.id);
    e.en[1] = get_edge_node(v1.id, v2.id);
    e.en[2] = get_edge_node(v2.id, v3.id);
    e.en[3] = get_edge_node(v3.id, v0.id);
}

void refine_quad(Element e)
{
    Node x0 = get_vertex_node(e.vn[0].id, e.vn[1].id);
    Node x1 = get_vertex_node(e.vn[1].id, e.vn[2].id);
    Node x2 = get_vertex_node(e.vn[2].id, e.vn[3].id);
    Node x3 = get_vertex_node(e.vn[3].id, e.vn[0].id);
}

```



```

e.sons[0] = create_quad(e.vn[0], x0, mid, x3);
e.sons[1] = create_quad(x0, e.vn[1], x1, mid);
e.sons[2] = create_quad(mid, x1, e.vn[2], x2);
e.sons[3] = create_quad(x3, mid, x2, e.vn[3]);
e.active = 0;
}

```

Each hash table is implemented as an array of linked lists. This hash table organization has the advantage of simple node removal, which is required if a node is no longer needed by any element.

3.3.3 Finding all neighbors of an element

In element-wise assembling procedure with arbitrary-level hanging nodes, when calculating fluxes across mesh edges, one needs values from both sides of the edge. To achieve this, the integration points need to be matched correctly, and proper values in these points from the adjacent element need to be extracted for all involved functions defined on the mesh (basis functions, previous time level solutions). This is not easy to manage. A simplified version of the algorithm is in Appendix 1.

3.3.4 Determining node type

Figure 3.13 shows a simple mesh with hanging nodes. Five types of nodes can be identified: standard vertex nodes (A), standard edge nodes (B), constrained vertex nodes (C), constrained edge nodes (D) and finally standard edge nodes and constrained vertex nodes at the same place (E).

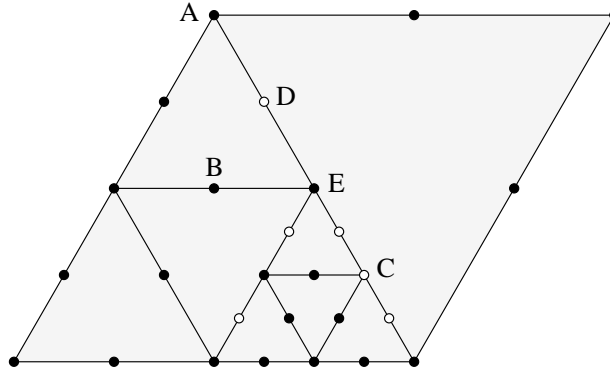


Figure 3.13: Mesh with hanging nodes

As the elements in the mesh get refined (or derefined), some nodes are created, others are removed and some change their types. One has to be able to quickly recognize whether a node is constrained or unconstrained without searching in the element hierarchy. This is achieved through the member `ref` of the structure `Node`. At all times this variable holds the number of elements pointing to the node. This

is the reason why all nodes of an element must be *referenced* (**ref** increased) when creating the element and *un-referenced* (**ref** decreased) when the element is being refined or deleted, as shown in subsection 3.3.2.

Now the cases (A) and (C) can be distinguished just by looking at **ref**. If **ref** = 3, the vertex node is constrained, if **ref** > 3 it is not constrained. Top-level vertex nodes have **ref** artificially increased to a large number, which ensures that they are always treated as unconstrained. Similarly, if **ref** = 1 for an edge node, it is constrained (D); if **ref** = 2 it is unconstrained (D).

The case (E) can be detected by calling the function `peek*_node`, which works the same way as `get*_node`, with the exception that the node is not created if it does not exist.

A node is deleted any time its **ref** reaches zero.

3.4 Automatic hp -adaptivity

The major difference between adaptivity in standard FEM and adaptivity in hp -FEM is the large number of element refinement options in the latter case. In standard h -adaptivity, elements can only be subdivided in space, e.g., using the *red-green* refinement technique (see [?] and the references therein). In hp -adaptivity, one can either increase the polynomial degree of elements without subdividing them in space or one can split elements in space and distribute the polynomial degree on the subelements in multiple ways.

3.4.1 Reference solution

The presence of many hp -refinement options per element implies that classical error estimates (in the form of one number per element) do not provide enough information to guide hp -adaptivity. Instead, one needs to use some information about the *shape* of the error function $e_{h,p} = u - u_{h,p}$. In principle, this information could be obtained by estimating higher derivatives, but such approach is not very practical. Usually it is more convenient to use a *reference solution*, i.e., an approximation u_{ref} which is at least one order more accurate than the coarse mesh solution $u_{h,p}$. The hp -adaptivity is then guided by an a-posteriori error estimate of the form

$$e_{h,p} = u_{ref} - u_{h,p} \quad (3.3)$$

Such approach, as opposed to classical error estimates, is virtually equation-independent. In this work we are constructing the enriched finite element space by uniform subdivision of all mesh elements and by increasing all element polynomial degrees by one, i.e.

$$u_{ref} = u_{h/2,p+1}$$

The reference solution u_{ref} can be obtained efficiently by utilizing information about lower frequencies from $u_{h,p}$ [?, ?, ?] if iterative solvers are used or by reconstructing the LU decomposition of the stiffness matrix from the LU-decomposed coarse mesh stiffness matrix when using direct sparse solvers, e.g. UMFPACK. However, we have not attempted any such optimization in this work, even though we admit that this issue must be addressed in order for the methods to be usable in practice.

3.4.2 Outline of the algorithm

The outer loop of our adaptivity algorithm is formally similar to the one of L. Demkowicz et al. (see [?]), with minor differences. The outline of the algorithm is as follows:

- A) Assume an initial coarse mesh $\mathcal{T}_{h,p}$ consisting of piecewise-constant, linear or quadratic elements. User input includes

- (a) prescribed relative tolerance $TOL > 0$ for the rate of energy norm of the approximate error function (3.3) and energy norm of the approximate solution
 - (b) threshold ERRT specifying how many elements should be refined in each hp -adaptivity step
- B) User input may also include
- (a) maximum increase of degrees of freedom MAX_STEP_NDOF in each adaptivity step
OR
 - (b) maximum number of degrees of freedom MAX_NDOF
- C) Create a fine mesh, $\mathcal{T}_{h/2,p+1}$
- D) Compute fine mesh approximation $u_{h/2,p+1} \in V_{h/2,p+1}$ on $\mathcal{T}_{h/2,p+1}$.
- E) Project the fine mesh approximation onto the coarse mesh $\mathcal{T}_{h,p}$. In this way, obtain a coarse mesh approximation.
- F) Calculate the desired norm (H^1 norm, H^1 seminorm, L^2 norm, custom norm) of the fine mesh approximations $NORM_i$ on every element K_i in the mesh, $i = 1, 2, \dots, M$. Construct the approximate error function (3.3) as the difference between the fine and coarse mesh approximations, calculate its energy norm ERR_i on every element K_i in the mesh, $i = 1, 2, \dots, M$. Calculate the global error estimate
- $$ERR = \left(\sum_{i=1}^M ERR_i \right) / \left(\sum_{i=1}^M NORM_i \right).$$
- G) If $ERR \leq TOL$, stop computation and proceed to postprocessing.
- H) Sort all elements into a list L according to their value $ERR_i/NORM_i$ in decreasing order.
- I) Determine the maximum of element errors, $ERR_{max} = \max\{ERR_i/NORM_i\}$, by taking the first item of L .
- J) Let $NDOF$ be the current number of degrees of freedom of $\mathcal{T}_{h,p}$. Repeat the following cycle:
- (a) If MAX_STEP_NDOF is set, and if the already added degrees of freedom in this hp -adaptivity step is greater than MAX_STEP_NDOF , break the cycle

- (b) Take the next element K_i from the list L
- (c) If $ERR_i/NORM_i < ERRT \cdot ERR_{max}$, break the cycle
- (d) Perform hp -refinement of K_i (to be described in more detail below).
- (e) If MAX_NDOF is set, and if the total number of degrees of freedom is now greater than MAX_NDOF , break the cycle

K) Continue with step C.

The crucial issue in the outer loop is determining how many elements should be refined. In [?] all elements whose error ERR is greater than 70% of the maximum error ERR_{max} are taken (ie. $ERRT = 0.7$). This is a relatively conservative choice which leads to good convergence curves, but may result in too many steps of the algorithm, which should be avoided as the evaluation of u_{ref} is very expensive. Our experience shows that setting $ERRT$ as low as 15% results in equally good convergence curves for most problems, but achieved in much fewer steps and in turn in much shorter time. Occasionally, however, the low value of $ERRT$ may cause too many elements to be refined and thus we introduced the limit MAX_STEP_NDOF of the number of degrees of freedom added in each step. Typically, one should not increase the number of degrees of freedom in each step by a large number.

3.4.3 Selecting optimal hp -refinement

The algorithm for determining optimal hp -refinements in [?] is quite complex. Built upon the projection-based interpolation theory, it first finds optimal hp -refinement of mesh edges using a 1D version of the hp -adaptive algorithm. The edge refinements then determine h -refinement of elements. Finally, best polynomial degrees for element interiors are selected. Each step of the algorithm represents a considerable implementation burden.

We have implemented a much simpler scheme, which is local in the sense that element refinements are selected without regard to the refinements of neighboring elements. For all elements $K \in \mathcal{T}_{h,p}$ of polynomial degree p picked by the outer loop we consider the following $N = 83$ hp -refinement options:

- A) Increase the polynomial degree of K to $p + 1$ without spatial subdivision.
- B) Increase the polynomial degree of K to $p + 2$ without spatial subdivision.
- C) Split K into four similar triangles K_1, K_2, K_3, K_4 . Define p_0 to be the integer part of $(p + 1)/2$. For each K_i , $1 \leq i \leq 4$ consider polynomial degrees $p_0 \leq p_i \leq p_0 + 2$.

For each of these N options we perform a standard e.g. L_2 -projection of the reference solution u_{ref} onto the corresponding polynomial space. The candidate with smallest projection error is selected for hp -refinement of the element K .

Each of the N projection problems requires the solution of a small system of linear algebraic equations with a symmetric positive-definite matrix. The solution of these systems can be further optimized by exploiting the incremental nature of the Cholesky decomposition algorithm and the fact that the spaces in item 3 above are partially embedded.

Example We again solve the model problem (2.1) with the same settings as in Paragraph 2.1. This time, automatic hp -adaptivity will be used. Although the maximum value of the exact solution is 1.0, due to the presence of undershoots and overshoots, the maximum values here are higher, as was seen already in section 2.4. This can be avoided by using a suitable shock capturing method, as it is done in the last chapter with numerical examples for the Euler equations.

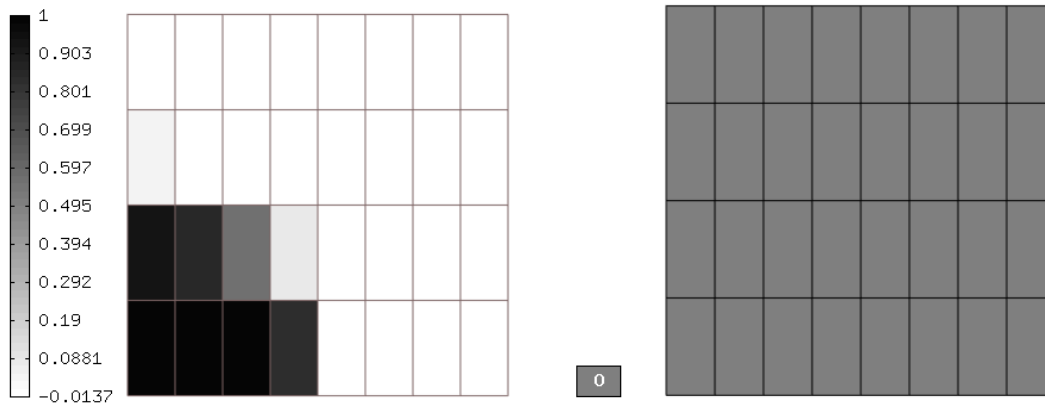


Figure 3.14: Solution and a mesh, relative error between the fine mesh and the coarse mesh approximation: 25.7746%.

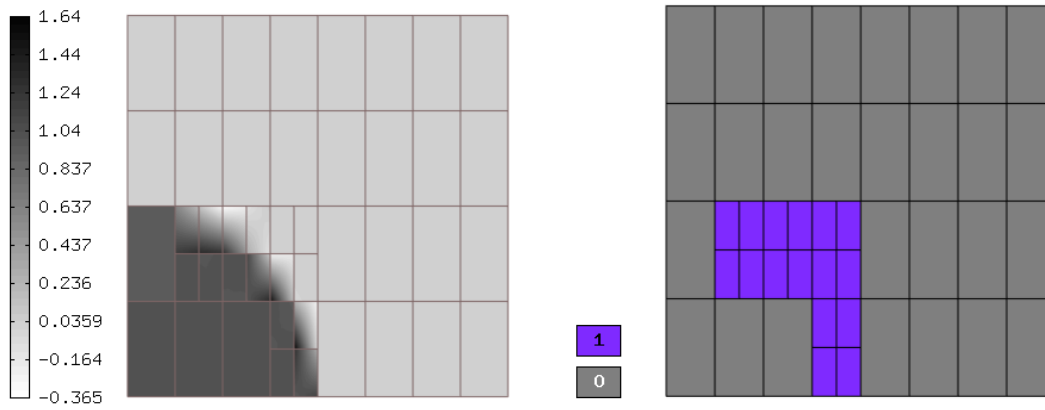


Figure 3.15: Solution and a mesh, relative error between the fine mesh and the coarse mesh approximation: 10.3565%.

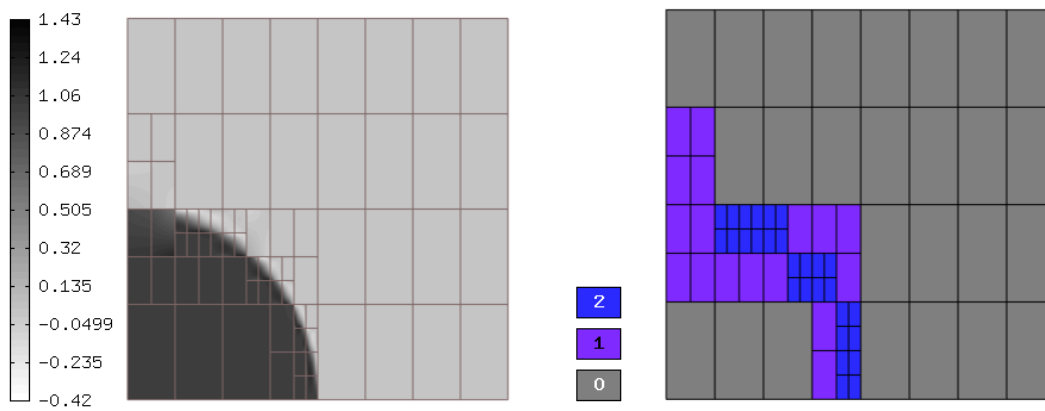


Figure 3.16: Solution and a mesh, relative error between the fine mesh and the coarse mesh approximation: 5.58551%.

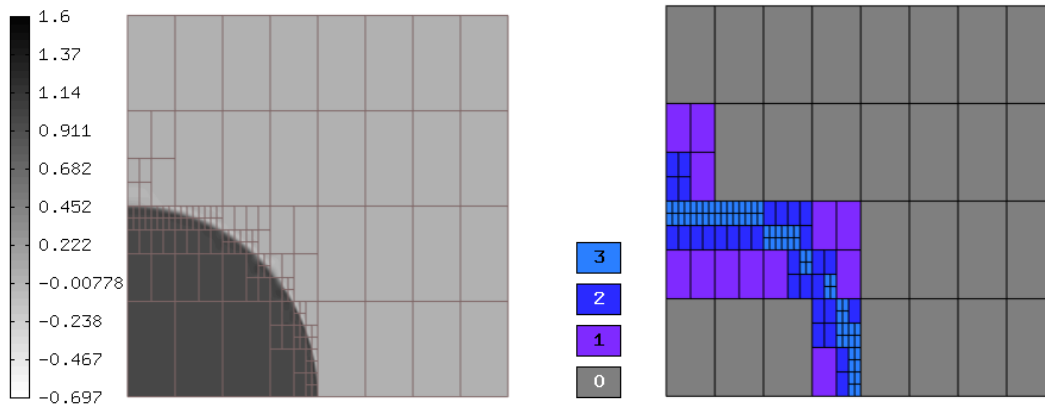


Figure 3.17: Solution and a mesh, relative error between the fine mesh and the coarse mesh approximation: 2.43584%.

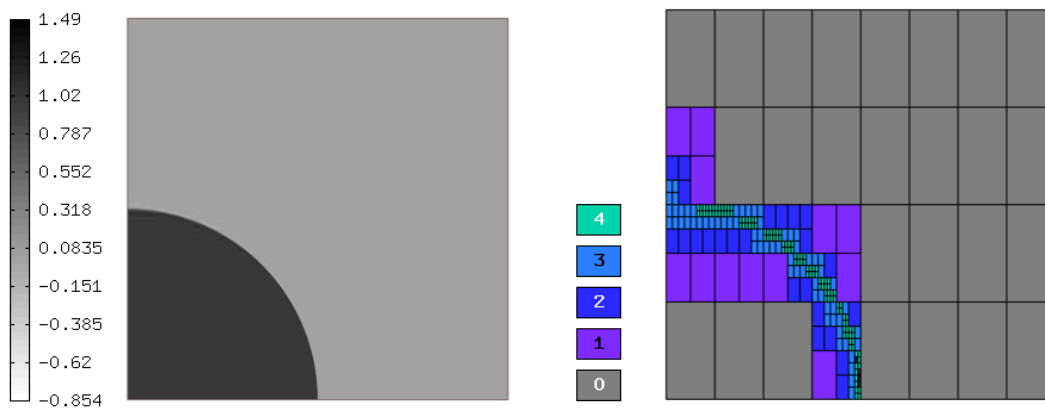


Figure 3.18: Solution and a mesh, relative error between the fine mesh and the coarse mesh approximation: 1.27544%.

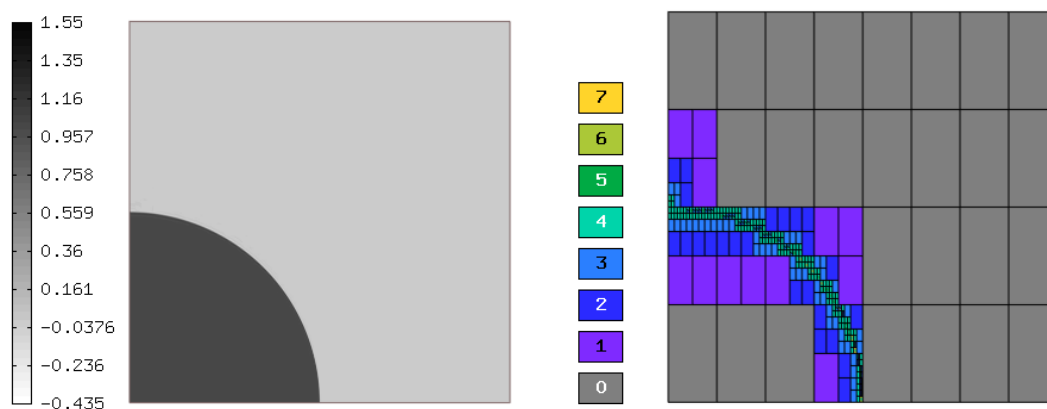


Figure 3.19: Solution and a mesh, relative error between the fine mesh and the coarse mesh approximation: 0.919277%.

3.4.4 Automatic *hp*-adaptivity for time-dependent problems

In the previous chapters we were concerned with the stationary problems. Let us now discuss automatic adaptivity for transient problems. Time-dependent coupled problems are challenging since one has to capture transient phenomena with sufficient accuracy, while the size of the problem must stay reasonably small. This leads to an obvious need for dynamically changing meshes for such problems. In most scientific computations, where dynamical meshes are used for time-dependent problems, data-transfer methods are necessary to move solution values between meshes for different time levels. As shown in subsection ??, data-transfer methods cause additional error and in case of time-dependent problems repeated transfers (usually simple interpolation) between meshes can have disastrous consequences.

We use orthogonal projections of the solution of transient problems using dynamical *hp*-meshes obtained fully automatically by the *hp*-adaptive algorithm. With the automatic *hp*-adaptivity and dynamical meshes the question of coarsening meshes between time levels arises. Mesh derefinement is particularly important in problems where sharp fronts (internal layers) move through the domain leaving smooth solutions behind them. We propose an original coarsening algorithm suitable for *hp*-FEM based on the *super-coarse solution*, which results in substantially fewer adaptive iterations.

The adaptive *hp*-FEM algorithm for time-dependent problems we use is obtained by combining the classical Rothe's method for time discretization with adaptive *hp*-FEM for the space discretization. The Rothe's method is a natural counterpart of the widely used Method of Lines (MOL). Recall that the MOL performs discretization in space while keeping the time variable continuous, which leads to a system of ODEs in time. The Rothe's method, on the contrary, preserves the continuity of the spatial variable while discretizing time. In every time step, an evolutionary PDE is approximated by means of one or more time-independent ones. For one step methods (such as implicit Runge-Kutta), the number of the time-independent equations per time step is proportional to the order of accuracy of the time discretization method. For example, when employing the implicit Euler method, one has to solve just one time-independent PDE per time step:

$$\frac{\partial u}{\partial t} = F(t, u) \quad \Rightarrow \quad \frac{u^{n+1} - u^n}{\Delta t} = F(t^{n+1}, u^{n+1}) \quad (3.4)$$

The Rothe's method is fully equivalent to the MOL if no adaptivity in space or time takes place, but it provides a better setting for the application of spatially adaptive algorithms. The spatial discretization error can be controlled by solving the time-independent equations adaptively, and the size of the time step can be adjusted using standard ODE techniques [?, ?, ?].

In this subsection, first the notion of *dynamical meshes* will be presented, then the refining and coarsening algorithms for time dependent problems are shown.

Dynamical meshes By \mathcal{T}_m let us denote a uniform coarse mesh covering the computational domain Ω . This mesh (called *master mesh*) is shared by the solution at all time levels, in other words all meshes can be obtained from this mesh by elementary refinements. At each time instant t_n an optimal mesh is found to suit the best the solution $\mathbf{u}^n(\mathbf{x})$. On the $(n+1)$ st time level, the approximated solution $\mathbf{u}^n(\mathbf{x})$, that has been obtained in the previous time step, is used as data. Note, however, that \mathbf{u}^n is defined on a locally refined mesh that was created automatically during the n th time step, while the unknown \mathbf{u}^{n+1} is solved adaptively starting from a coarser mesh. As a result, the meshes obtained on each time level are different, i.e., the mesh changes dynamically in time.

As mentioned above, in transient problems the optimal meshes are on each time level obtained automatically and they can change from one time step to another, which induces need for both refining and coarsening strategies.

Refining algorithm

In Chapter ?? the goal of adaptive algorithm was to decrease the error as low as possible. On the other hand for time-dependent problems we want to sustain the space error on approximately the same level, which would result into meshes with smoothly changing number of degrees of freedom. For time-dependent problems, as well as for time-independent, the amount of elements refined in one adaptive iteration will depend on the global solution error estimate in that iteration. As described, one can drive the refining algorithm by variables TOL , $ERRT$, MAX_STEP_NDOF , MAX_NDOF .

Coarsening algorithm

The most primitive strategy to coarsen the mesh on the next time level is to start on each level from the very coarse (master) mesh and perform the automatic adaptivity to find the optimal mesh for a particular solution $\mathbf{u}^n(\mathbf{x})$. Although this would result in the most optimal meshes in each iteration, it is virtually impossible to carry it out due to the immense computational time demands. Moreover, a lot of work would be wasted since the solutions from two adjacent time steps usually differ only mildly even though the solution changes significantly in the whole time domain. Global derefinement would result in similar difficulties.

We present a coarsening algorithm that prepares the mesh for the adaptive algorithm on the next time level by local derefinements. Thus, we remove only unnecessary refinements from previous time levels. In this way the meshes for particular time steps are suboptimal, but the number of adaptive iterations performed in each time step significantly decreases.

In higher-order finite element method we have two questions. First question is which elements can be coarsened and second is how to coarsen them. In *hp*-FEM we can either decrease the polynomial order of the element or derefine the element

geometrically. We can also employ various choices for assigned polynomial order. Situation is depicted in Fig. 3.20.

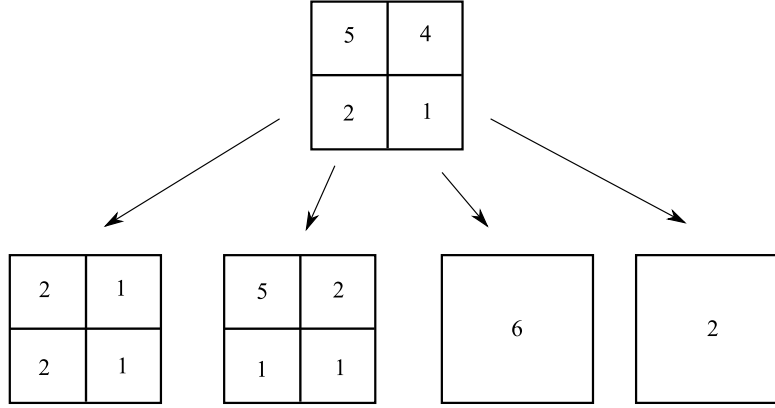


Figure 3.20: Coarsening choices for an element.

However, the coarsening serves only to prepare the mesh for the next adaptive process, thus we do not have to search for the best unrefinement as we do in case of refinements. It is perfectly sufficient to perform any coarsening that does not cause error increase and the mesh will optimize itself in the next adaptive loop. So far, to keep the algorithm reasonably simple we allow only coarsening as reversing of previous element refinements. In this way the tree-like structure of meshes starting from the *master mesh* is preserved. Unfortunately, the coarsening algorithm cannot be just an opposite to the refining. While with the refining algorithm we are seeking for the best candidate in the sense of error decrease, here we are looking for a candidate whose error will not exceed some tolerance (so that it would not be subsequently refined).

Let us denote by e_{max} the error of the element with maximal error from the last adaptive iteration. From Alg. ?? from Chapter ?? we know that errors of all elements are below ERR_{max} and that it should be satisfied as well after the coarsening procedure. First we calculate so called *super-coarse solution* for polynomial orders. By that we mean that polynomial order on all elements is decreased by one. We calculate error estimates for all elements in the super-coarse mesh with respect to the reference solution and lower the polynomial order on those elements whose error after coarsening is less than $k * ERR_{max}$, where parameter $k \in (0, 1)$ is chosen to ensure that element will not be subsequently refined in the adaptive procedure on the next time level. Similar procedure is run for a spatial coarsening – *super-coarse solution* is calculated on the mesh where all elements with four active subelements are coarsened and polynomial order on such elements is assigned to be the maximum of orders on four subelements. In a similar way as before we determine which elements can be also spatially unrefined without significant increase of the error.

The whole procedure for time-dependent problems with the refining and coarsening strategy is described in Alg. 1. Effectivity of the approach is demonstrated on the numerical examples in the last chapter of this thesis.

Algorithm 1: Adaptive algorithm for time-dependent problems with improved stopping criterion.

```

foreach time level do
  repeat
    compute reference solution  $u_{ref}$ ;
    project to get a solution on current mesh  $u$ ;
    evaluate global error estimate  $ERR_{gl}$  and local errors on elements  $ERR_i$ ;
    if  $ERR_{gl} < TOL$  then
      done = true;
      break;
    else
      sort all elements by their error;
      processed error = 0;
      foreach element do
        if ( $processed\ error > c * (ERR_{gl} - TOL)$ ) or
          ( $ERR_i < k * ERR_{max}$ ) then
          done = true;
          break;
        else
          find optimal refinement;
          processed error +=  $ERR_i$ ;
      until not done;
    calculate super-coarse solution for polynomial orders;
    decrease polynomial orders when possible;
    calculate super-coarse solution for spatial refinements;
    geometrically coarsen elements when possible;

```

Chapter 4

Dicontinuous Galerkin discretization of the Euler equations

For what follows, we introduce the following notation. By $\varphi|_{\Gamma_{ij}}$ and $\varphi|_{\Gamma_{ji}}$ we denote the values of φ on Γ_{ij} considered from the interior and the exterior of K_i respectively. The symbols

$$\langle \varphi \rangle_{ij} = \frac{1}{2} (\varphi|_{\Gamma_{ij}} + \varphi|_{\Gamma_{ji}}), \quad [\varphi]_{ij} = \varphi|_{\Gamma_{ij}} - \varphi|_{\Gamma_{ji}} \quad (4.1)$$

stand for the average of the two values $\varphi|_{\Gamma_{ij}}$ and $\varphi|_{\Gamma_{ji}}$, and their difference (i.e. jump of φ on Γ_{ij}) respectively.

4.1 Numerical Fluxes

The DGFEM discretization of the Euler equations in their conservative form follows the principles of the section 2.4. We take the equations in the form (1.73). We leave the time derivative for now, and we need to find suitable numerical fluxes to approximate the fluxes $\partial \mathbf{f}_x, \partial \mathbf{f}_y$ through the faces Γ_{ij} . The numerical fluxes will be sought so that the requirements from the subsection 2.2.2 are met.

Construction of some numerical fluxes The following type of numerical fluxes are usually called *flux vector splitting schemes* and we use the knowledge from the section 1.2, namely (1.81), (1.83), ..., and (1.102). On the basis of (1.102) we define the matrices

$$\Lambda^\pm = \text{diag}(\lambda_1^\pm, \dots, \lambda_m^\pm), \quad |\Lambda| = \text{diag}(|\lambda_1|, \dots, |\lambda_m|). \quad (4.2)$$

Also we define

$$\mathbb{P}^\pm = \mathbb{P} \Lambda^\pm \mathbb{P}^{-1}, \quad |\mathbb{P}| = \mathbb{P} |\Lambda| \mathbb{P}^{-1}. \quad (4.3)$$

These matrices depend on $\mathbf{w} \in D$ and $\mathbf{n} \in \mathcal{S}_1$. Now we define the following two schemes that will later be used in the numerical examples.

A) *The Steger-Warming scheme* has the numerical flux:

$$\mathbf{H}_{SW}(\mathbf{u}, \mathbf{v}, \mathbf{n}) = \mathbb{P}^+(\mathbf{u}, \mathbf{n}) \mathbf{u} + \mathbb{P}^-(\mathbf{v}, \mathbf{n}) \mathbf{v}, \quad \mathbf{u}, \mathbf{v} \in D, \quad \mathbf{n} \in \mathcal{S}_1.$$

This scheme is rather diffusive, so another scheme is preferred.

B) *The Vijayasundaram scheme* has the numerical flux:

$$\mathbf{H}_V(\mathbf{u}, \mathbf{v}, \mathbf{n}) = \mathbb{P}^+\left(\frac{\mathbf{u} + \mathbf{v}}{2}, \mathbf{n}\right) \mathbf{u} + \mathbb{P}^-\left(\frac{\mathbf{u} + \mathbf{v}}{2}, \mathbf{n}\right) \mathbf{v}, \quad \mathbf{u}, \mathbf{v} \in D, \quad \mathbf{n} \in \mathcal{S}_1.$$

This scheme contains a *partial upwinding*, where the flux is computed at the point $x_{i+\frac{1}{2}}$ ($x_{i-\frac{1}{2}}$) with the use of the value \mathbf{w}_i^k or \mathbf{w}_{i+1}^k (\mathbf{w}_{i-1}^k or \mathbf{w}_i^k) corresponding to the mesh point located in the upwind direction wrt. the propagation speed given by the eigenvalues λ_i . In the Steger-Warming scheme we speak of a *full upwinding*.

Another possible construction of the numerical flux is based on the rotational invariance of the Euler equations, see Section 1.2.3. We introduce a new Cartesian coordinate system $\tilde{x}_1, \dots, \tilde{x}_N$ with the origin at the midpoint of the face Γ , the coordinate \tilde{x}_1 oriented in the direction of the normal \mathbf{n} and $\tilde{x}_2, \dots, \tilde{x}_N$ tangent to Γ . The rotational invariance of the Euler equations implies that these equations have the form

$$\frac{\partial \mathbf{q}}{\partial t} + \sum_{s=1}^2 \frac{\partial \mathbf{f}_s(\mathbf{q})}{\partial \tilde{x}_s} = 0, \quad (4.4)$$

where

$$\mathbf{q} = \mathbb{Q}(\mathbf{n}) \mathbf{w}$$

with matrix $\mathbb{Q}(\mathbf{n})$ defined in (1.92). The tangential derivatives are neglected, and only the system with one space variable is considered:

$$\frac{\partial \mathbf{q}}{\partial t} + \frac{\partial \mathbf{f}_1(\mathbf{q})}{\partial \tilde{x}_1} = 0. \quad (4.5)$$

Then the numerical flux is defined in the form

$$\mathbf{H}(\mathbf{u}, \mathbf{v}, \mathbf{n}) = \mathbb{Q}^{-1}(\mathbf{n}) \mathbf{g}_R(\mathbb{Q}(\mathbf{n}) \mathbf{u}, \mathbb{Q}(\mathbf{n}) \mathbf{v}), \quad (4.6)$$

where $\mathbf{g}_R = \mathbf{g}_R(\mathbf{q}_1, \mathbf{q}_2)$ is a numerical flux (so-called *approximate Riemann solver*, see [?]) for the system (4.5) with one space variable. This numerical flux is complicated to be linearized, and because we wanted to employ a semi-implicit method for time integration, we opted for the Vijayasundaram numerical flux in most cases.

4.2 Time discretization

This section closely follows [?]. Using numerical fluxes we can introduce new forms

$$\begin{aligned}
 (\mathbf{w}_h, \boldsymbol{\varphi}_h) &= \int_{\Omega} \mathbf{w}_h(\mathbf{x}) \cdot \boldsymbol{\varphi}_h(\mathbf{x}) d\mathbf{x}, \\
 \tilde{b}_h(\mathbf{w}_h, \boldsymbol{\varphi}_h) &= - \sum_{K \in \mathcal{T}_h} \int_K \sum_{s=1}^2 \mathbf{f}_s(\mathbf{w}_h(\mathbf{x})) \cdot \frac{\partial \boldsymbol{\varphi}_h(\mathbf{x})}{\partial x_s} d\mathbf{x} \\
 &\quad + \sum_{K_i \in \mathcal{T}_h} \sum_{j \in \mathcal{S}(i)} \int_{\Gamma_{ij}} \mathbf{H}(\mathbf{w}(t)|_{\Gamma_{ij}}, \mathbf{w}(\mathbf{x})|_{\Gamma_{ji}}, \mathbf{n}_{ij}) \cdot \boldsymbol{\varphi}_h(\mathbf{x}) dS
 \end{aligned} \tag{4.7}$$

for $\mathbf{w}_h, \boldsymbol{\varphi}_h \in [V_h]^4$. We say that \mathbf{w}_h is the approximate solution of (1.73) in $Q_T = \Omega \times (0, T)$, if it satisfies the conditions

$$\begin{aligned}
 \text{A)} \quad & \mathbf{w}_h \in \mathcal{C}^1([0, T], [V_h]^4)
 \end{aligned} \tag{4.8}$$

$$\begin{aligned}
 \text{B)} \quad & \frac{d}{dt}(\mathbf{w}(t), \boldsymbol{\varphi}_h) + \tilde{b}_h(\mathbf{w}_h(t), \boldsymbol{\varphi}_h) = 0 \quad \forall \boldsymbol{\varphi}_h \in [V_h]^4 \forall t \in (0, T)
 \end{aligned} \tag{4.9}$$

$$\begin{aligned}
 \text{C)} \quad & \mathbf{w}_h(0) = \Pi_h \mathbf{w}^0
 \end{aligned} \tag{4.10}$$

where $\Pi_h \mathbf{w}^0$ is the L^2 -projeftion of the initial condition \mathbf{w}^0 onto the space $[V_h]^4$. If we set $p = 0$, then we obviously obtain the finite volume method.

Relations (4.9) represent a system of ordinary differential equations which can be solved by a suitable numerical method. Since we are interested in applying the Rothe's method, we now want to discretize the time derivative. In order to do so, we consider a partition $0 = t_0 < t_1 < t_2 < \dots$ of the time interval $(0, T)$ and set $\tau_k = t_{k+1} - t_k$. We use the notation \mathbf{w}_h^k for the approximation of $\mathbf{w}_h(t_k)$. Then we apply the simple implicit *backward Euler method* and our *discrete problem* reads: for each $k \geq 0$ find \mathbf{w}_h^{k+1} such that

$$\begin{aligned}
 \text{A)} \quad & \mathbf{w}_h^{k+1} \in [V_h]^4
 \end{aligned} \tag{4.11}$$

B)

$$\left(\frac{\mathbf{w}_h^{k+1} - \mathbf{w}_h^k}{\tau_k}, \boldsymbol{\varphi}_h \right) + \tilde{b}_h(\mathbf{w}_h^{k+1}, \boldsymbol{\varphi}_h) = 0 \quad \forall \boldsymbol{\varphi}_h \in [V_h]^4, \quad k = 0, 1, \dots \quad (4.12)$$

C)

$$\mathbf{w}_h^0 = \Pi_h \mathbf{w}^0 \quad (4.13)$$

This scheme leads to a system of highly nonlinear algebraic equations whose numerical solution is rather complicated. In order to simplify the problem, in the following we shall linearize relations (4.12) and obtain a linear system.

4.2.1 Linearization

By (4.7), for $\mathbf{w}_h^{k+1}, \boldsymbol{\varphi}_h \in [V_h]^4$ we have

$$\tilde{b}_h(\mathbf{w}_h^{k+1}, \boldsymbol{\varphi}_h) = - \sum_{K \in \mathcal{T}_h} \int_K \sum_{s=1}^2 \mathbf{f}_s(\mathbf{w}_h^{k+1}(\mathbf{x})) \cdot \frac{\partial \boldsymbol{\varphi}(\mathbf{x})}{\partial x_s} d\mathbf{x} \quad (=:\tilde{\sigma}_1) \quad (4.14)$$

$$+ \sum_{K_i \in \mathcal{T}_h} \sum_{j \in S(i)} \int_{\Gamma_{ij}} \mathbf{H}(\mathbf{w}_h^{k+1}|_{\Gamma_{ij}}, \mathbf{w}_h^{k+1}|_{\Gamma_{ji}}, \mathbf{n}_{ij}) \cdot \boldsymbol{\varphi}_h(\mathbf{x}) dS \quad (=:\tilde{\sigma}_2) \quad (4.15)$$

The terms $\tilde{\sigma}_1$, and $\tilde{\sigma}_2$ are linearized as follows. We set

$$\tilde{\sigma}_1 \approx \sigma_1 = \sum_{K \in \mathcal{T}_h} \int_K \sum_{s=1}^2 \mathbb{A}_s(\mathbf{w}_h^k(\mathbf{x})) \mathbf{w}_h^{k+1}(\mathbf{x}) \cdot \frac{\partial \boldsymbol{\varphi}_h(\mathbf{x})}{\partial x_s} d\mathbf{x} \quad (4.16)$$

and

$$\tilde{\sigma}_2 \approx \sigma_2 = \sum_{K_i \in \mathcal{T}_h} \sum_{j \in S(i)} \int_{\Gamma_{ij}} \left[\mathcal{P}^+ \left(\langle \mathbf{w}_h^k \rangle_{ij}, \mathbf{n}_{ij} \right) \mathbf{w}_h^{k+1}|_{\Gamma_{ij}} \right. \quad (4.17)$$

$$\left. + \mathcal{P}^- \left(\langle \mathbf{w}_h^k \rangle_{ij}, \mathbf{n}_{ij} \right) \mathbf{w}_h^{k+1}|_{\Gamma_{ji}} \right] \cdot \boldsymbol{\varphi}_h dS \quad (4.18)$$

$$+ \sum_{K_i \in \mathcal{T}_h} \sum_{j \in S(i) \setminus s(i)} \int_{\Gamma_{ij}} \mathbf{H}(\mathbf{w}_h^{k+1}|_{\Gamma_{ij}}, \mathbf{w}_h^{k+1}|_{\Gamma_{ji}}, \mathbf{n}_{ij}) \cdot \boldsymbol{\varphi}_h(\mathbf{x}) dS. \quad (4.19)$$

We omit the linearization of the fluxes across domain boundaries, for details see [?]. Finally, we define the form

$$b_h(\mathbf{w}_h^k, \mathbf{w}_h^{k+1}, \boldsymbol{\varphi}_h) = -\sigma_1 + \sigma_2. \quad (4.20)$$

The form is linear wrt. the second and third variable. Using this form we come to the following *semi-implicit linearized numerical scheme*: for each $k \geq 0$ find \mathbf{w}_h^{k+1} such that

$$\text{A)} \quad \mathbf{w}_h^{k+1} \in [V_h]^4 \quad (4.21)$$

$$\text{B)} \quad (\mathbf{w}_h^{k+1}, \boldsymbol{\varphi}_h) + \tau_k b_h(\mathbf{w}_h^k, \mathbf{w}_h^{k+1}, \boldsymbol{\varphi}_h) = (\mathbf{w}_h^k, \boldsymbol{\varphi}_h) \quad \forall \boldsymbol{\varphi}_h \in [V_h]^4, \quad k = 0, 1, \dots \quad (4.22)$$

$$\text{C)} \quad \mathbf{w}_h^0 = \Pi_h \mathbf{w}^0 \quad (4.23)$$

4.3 Boundary conditions

In all examples in the next chapter, the boundary of the computational domain Ω is divided into three parts.

A) On a fixed impermeable wall $\Gamma_W \subset \partial\Omega$ we use the condition $\mathbf{v} \cdot \mathbf{n} = 0$. Then the flux $\mathcal{P}(\mathbf{w}, \mathbf{n})$ has the form

$$\mathcal{P}(\mathbf{w}, \mathbf{n}) = \sum_{s=1}^2 \mathbf{f}_s(\mathbf{w}) n_s \quad (4.24)$$

$$= (\mathbf{v} \cdot \mathbf{n}) \mathbf{w} + p(0, n_1, n_2, \mathbf{v} \cdot \mathbf{n})^T \quad (4.25)$$

$$= p(0, n_1, n_2, 0)^T, \quad (4.26)$$

which is uniquely determined on Γ_W by the extrapolated value of the pressure, i.e. by $p_j^k := p_i^k$. Therefore, on the part Γ_W of the boundary we define the numerical flux

$$\mathbf{H}(\mathbf{w}_i^k, \mathbf{w}_j^k, \mathbf{n}) = p_i^k(0, n_1, n_2, 0)^T. \quad (4.27)$$

We can see that on the impermeable part of the boundary, 2 eigenvalues λ_2, λ_3 are zero, and the eigenvalue λ_1 is negative, and the eigenvalue λ_4 is positive. We prescribe only $\mathbf{v} \cdot \mathbf{n} = 0$, and extrapolate the pressure.

B) On the inlet

$$\Gamma_I(t) = \{x \in \partial\Omega; \mathbf{v}(x, t) \cdot \mathbf{n}(x) < 0\}$$

and the outlet

$$\Gamma_O(t) = \{x \in \partial\Omega; \mathbf{v}(x, t) \cdot \mathbf{n}(x) > 0\}$$

parts of the boundary it is necessary to use nonreflecting boundary conditions. In this case we used the so-called *characteristics-based* boundary

conditions, according to [?]. Using the rotational invariance, we tranform the Euler equations to the coordinates \tilde{x}_1 , in the direction of the outer normal \mathbf{n} to the boundary, and \tilde{x}_2 , tangential to the boundary and linearize the resulting system around the state $\mathbf{q}_{ij} = \mathbb{Q}(\mathbf{n}_{ij})$. Then we obtain the linear system

$$\frac{\partial \mathbf{q}}{\partial t} + \mathbb{A}_1(\mathbf{q}_{ij}) \frac{\partial \mathbf{q}}{\partial \tilde{x}_1} = 0, \quad (4.28)$$

for the vector-valued function $\mathbf{q} = \mathbb{Q}(\mathbf{n}_{ij}) \mathbf{w}$, considered in the set $(-\infty, 0) \times (0, \infty)$ and equipped with the initial and boundary conditions

$$\mathbf{q}(\tilde{x}_1, 0) = \mathbf{q}_{ij}, \quad \tilde{x}_1 < 0, \quad (4.29)$$

$$\mathbf{q}(0, t) = \mathbf{q}_{ij}, \quad t > 0. \quad (4.30)$$

The goal is to choose \mathbf{q}_{ij} in such a way that this initial-boundary value problem is well-posed, i.e. has a unique solutions. The method of characteristics leads to the following process:

Let us put $\mathbf{q}_{ji}^* = \mathbb{Q}(\mathbf{n}_{ij}) \mathbf{w}_{ji}^*$, where \mathbf{w}_{ji} is a prescribed boundary state at the inlet or outlet. We calculate eigenvectors \mathbf{r}_s corresponding to the eigenvalues $\lambda_1, \dots, \lambda_4$, of the matrix $\mathbb{A}_1(\mathbf{q}_{ij})$, arrange them as columns in the matrix \mathbb{T} and calculate \mathbb{T}^{-1} (explicit formulae can be found in [?], section 3.1). Now we set

$$\alpha = \mathbb{T}^{-1} \mathbf{q}_{ij}, \quad \beta = \mathbb{T}^{-1} \mathbf{q}_{ji}^*. \quad (4.31)$$

and define the state \mathbf{q}_{ji} by the relations

$$\mathbf{q}_{ji} := \sum_{s=1}^4 \gamma_s \mathbf{r}_s, \quad \gamma_s = \begin{cases} \alpha_s, & \lambda_s \geq 0, \\ \beta_s, & \lambda_s < 0. \end{cases} \quad (4.32)$$

Finally, the sought boundary state \mathbf{w}_{ji} is defined as

$$\mathbf{w}|_{\Gamma_{ji}} = \mathbf{w}_{ji} = \mathbb{Q}^{-1}(\mathbf{n}_{ij}) \mathbf{q}_{ji}. \quad (4.33)$$

4.3.1 Semi-implicit method for the time integration

4.3.2 Stability and the CFL condition

Modry Feist, stranky viz TODOdiplomka na google docs

4.4 Shock capturing