

# Robust Algebraic Multigrid Methods in Magnetic Shielding Problems

## Diplomarbeit

zur Erlangung des akademischen Grades  
Diplom - Ingenieur  
in der Studienrichtung Technische Mathematik

eingereicht von  
Stefan Reitzinger

angefertigt am Institut für Analysis und Numerik  
der Technisch - Naturwissenschaftlichen Fakultät  
der Johannes Kepler Universität Linz

eingereicht bei  
O.Univ.-Prof. Dr. Ulrich Langer

Linz, Mai 1998

Für meine lieben Eltern

## Preface

During my study at the Johannes Kepler University in Linz I passed a project seminar, where I and a colleague of mine had to work out the simulation of magnetic shielding problems in 2D. In this work we discussed an appropriate mathematical model and presented some numerical results. Unfortunately the numerical calculations were not done very sophisticatedly and therefore only a unrealistic case could be calculated which hardly appears in real life. Nevertheless the mathematical model to describe a shielding phenomenon was well suited and actually only a robust solution technique was left.

One year later, after my study in Denmark, I took up the problem again and found out that a lot of practical problems lead to the same problem class. This motivated me to write a diplom thesis about robust solution strategies for magnetic shielding problems in 2D.

At this point I would like to express my thanks to Prof.Dr. Ulrich Langer who gave me the chance of writing this diplom thesis and for supporting me during this time. Additionally I express my thanks to his coworkers at the Institute of Analysis and Numerics at the Johannes Kepler University in Linz. Special thanks to Dr. Gundolf Haase and Dipl.Ing. Joachim Schöberl who gave me a lot of ideas and practical hints for implementation. Apart from that they prove great patience by answering my questions during the development phase of this thesis.

Last but not least I would like to express my thanks to Dipl.Ing. Wolfram Mühlhuber and Karin Swoboda for their help.

Mai 1998,

Stefan Reitzinger

## Abstract

The aim of this diplom thesis is to provide a robust and efficient solver for large sparse and poor conditioned linear systems arising from the FE-method for elliptic scalar PDEs of second order.

For a counter example the problem of magnetic shielding is used. Therefore the Maxwell's equations for stationary objects are reduced to a scalar PDE of second order with appropriate boundary conditions.

In order to solve the equation by means of FEM, a discretization for micro scales is introduced. Especially long thin elements are suggested to keep the number of unknowns small in areas of micro structures. Constructively a finite element analysis is carried out where also a convergence result of the FE-solution in  $H^1$  is presented.

To achieve an efficient and robust solution strategy the algebraic multigrid method of Ruge and Stüben is introduced. Additionally three different areas of application are presented for this AMG method, i.e. preconditioner, coarse grid solver for a full multigrid method, and black box solver.

Because this AMG method normally works well for M-matrices, a technique is presented to attain M-matrices, if the underlying linear system arises from an FE-discretization. The method to achieve the M-matrix property is based on the element matrices.

The algorithm was implemented as black box solver in the finite element package FEPP. Therein AMG was applied as preconditioner for the conjugate gradient method.

Some numerical experiments are presented, where long thin quadrilaterals are used with ratio of the longest and shortest side of 1 to  $10^{-3}$ . Additionally parameter jumps of order  $10^{-6}$  to  $10^{+6}$  are considered.

Concluding AMG has been proven, at least in a numerical way, to be an efficient and robust solver for magnetic shielding problems, if it is used as a preconditioner for the CG-method. If long thin quadrilaterals are used for discretization the modified preconditioner also behaves very robust.

## Zusammenfassung

Das Ziel dieser Diplomarbeit ist die Entwicklung einer robusten und effizienten Lösungsstrategie für große, dünnbesetzte und schlecht konditionierte lineare Gleichungssysteme, welche bei der Diskretisierung von elliptischen Differentialgleichungen 2. Ordnung entstehen.

Als Beispiele dienen Abschirmprobleme in 2D, wofür die Maxwell'schen Gleichungen als Grundlage für die Modellierung benutzt werden. Diese Gleichungen reduzieren sich zu einer skalaren partiellen Differentialgleichung 2. Ordnung mit geeigneten Randbedingungen.

Um diese Problemklasse numerisch mit Hilfe der Methode der Finiten Elemente zu lösen, wird eine geeignete Diskretisierung für Mikroskalen eingeführt. Zu diesem Zweck werden lange dünne Elemente vorgeschlagen, um die Anzahl der Unbekannten gering zu halten. Weiters wird eine Konvergenzanalyse für die FEM-Lösung in  $H^1$  bereitgestellt.

Als Grundlage für einen effizienten und robusten Löser dient die algebraische Multigrid Methode von Ruge und Stüben. Aufbauend auf dieser Methode werden drei Anwendungsgebiete beschrieben: Vorkonditionierer, Grobgitterlöser für die Full Multigrid Methode und Black Box Löser. Da die vorgeschlagene Technik im allgemeinen nur für M-Matrizen zufriedenstellend funktioniert wird eine Methode bereitgestellt, welche spektraläquivalente M-Matrizen generiert. Diese Technik basiert auf den Elementmatrizen.

Der Algorithmus wurde im Finite Elemente Programm FEPP implementiert, in dem AMG als Vorkonditionierer für das konjugierte Gradientenverfahren verwendet wird.

Abschließend werden numerische Experimente durchgeführt, in denen lange, dünne Vierecke mit einem Seitenverhältnis von ca.  $1 : 10^{-3}$  verwendet werden. Zusätzlich werden Parametersprünge in der Größenordnung von  $10^{-6}$  bis  $10^{+6}$  angenommen.

AMG ist ein, im zumindest numerischen Sinn, effizienter und robuster Löser für Abschirmprobleme, falls es als Vorkonditionierer für das CG Verfahren verwendet wird. Werden lange dünne Vierecke zur Diskretisierung verwendet, so ist der modifizierte Vorkonditionierer ebenfalls sehr robust.

# Contents

<b>1</b>	<b>Introduction</b>	<b>6</b>
<b>2</b>	<b>Problem description, formulation and preliminaries</b>	<b>9</b>
2.1	A practical problem and its motivation . . . . .	9
2.2	The physical problem . . . . .	10
2.3	The mathematical problem . . . . .	12
2.3.1	The scalar potential . . . . .	12
2.3.2	The vector potential . . . . .	14
2.4	The resulting variational form . . . . .	15
2.4.1	A weak formulation in 2D . . . . .	15
2.4.2	Existence and uniqueness . . . . .	16
2.4.3	Weak formulation for the vector-potential . . . . .	17
<b>3</b>	<b>Discretization in 2D</b>	<b>18</b>
3.1	General remarks . . . . .	18
3.2	The problem of micro scales . . . . .	19
3.3	Refinement strategy . . . . .	20
<b>4</b>	<b>Finite element analysis</b>	<b>22</b>
4.1	The resulting discrete problem - general remarks . . . . .	22
4.1.1	The finite element space . . . . .	22
4.1.2	The condition number . . . . .	23
4.2	The approximation error and a convergence result . . . . .	24
4.2.1	The approximation error . . . . .	24
4.2.2	A convergence result in $H^1$ . . . . .	25
<b>5</b>	<b>Numerical solution</b>	<b>27</b>
5.1	Motivation and practical usage of RSAMG . . . . .	27
5.1.1	RSAMG as a black box solver . . . . .	27
5.1.2	Iterative solver for FMG . . . . .	28
5.1.3	RSAMG as a preconditioner . . . . .	30
5.2	AMG of Ruge and Stüben (RSAMG) . . . . .	32
5.2.1	Terminology, assumptions and notations . . . . .	33
5.2.2	General convergence theory for the V-cycle . . . . .	35
5.2.3	The smoothing property . . . . .	36
5.2.4	The approximation property . . . . .	39
5.2.5	Interpolation formulas . . . . .	41
5.2.6	The coarsening algorithm and remarks for implementation . . . . .	42
5.2.7	The algorithm for general s.p.d. matrices . . . . .	43
<b>6</b>	<b>Numerical results</b>	<b>51</b>
6.1	First example . . . . .	51
6.2	Shielding problem with one thin material . . . . .	52
6.3	Shielding problem with different material layers . . . . .	55
<b>7</b>	<b>Conclusions and further remarks</b>	<b>58</b>

# 1 Introduction

In recent years the simulation of technical processes have become very trendy not only in the academic field but also industry realized the importance and effectivity of numerical simulation tools. The reason therefore is on one hand the permanent growing in performance and memory resources of computers allowing the computation of large problems. On the other hand there exists an advanced mathematical theory which provides fast and accurate solution strategies.

One branch of industry where this tendency can be observed is the electric and magnetic field computation. Of course there, and likewise in every branch of industry the complexity of the arising real-life problems exceeds often the feasibilities of the computers. Nevertheless also downsized problems are very helpful and important in a lot of cases. One of these problems in electric and magnetic field computations is the so called shielding problem, which appears very often in practice. Let us describe some of these problems to have a look at the complexity of this problem class.

- A data-line, which can be found in every antenna or phone cable can not be a normal cable. Because if it would be a normal cable the data sent via this cable could be strongly influenced and maybe disturbed by another electric field. Such fields arise for example in the neighborhood of every electric cable or simply from the earth. For these reasons a data-line has to be shielded.
- Another field of application is the shielding of sensitive measurements. For example in medicine we want to save the ECG from varying magnetic fields, because such fields could influence the measurements with dramatic consequences for the diagnosis. Another example is due to electric control devices. They may be very sensitive against an electric field and therefore have to be shielded.
- Apart from the examples above one could also want to produce an electric field in a certain area. This again could be in medicine to head up a part of a body.

The list of possible applications in shield computation and of course in electric field computation is much longer. However, all the problems mentioned above have to be downsized from the real-life problem to some reasonable problem size. This is a part of ‘Mathematical Modeling’ and can be found for example in [32]. For down sized mathematical models an analysis can be done and therefore several parameter settings can be carried out very easily. This is a great advantage compared to the classical development process, where for nearly every parameter setting a new prototype has to be built. With numerical simulation we can skip this phase by making experiments on the computer. Apart from the easy handling we save time and money. Straightforward such a development phase on the computer could be worked out to get an optimal parameter setting. In literature this is often referred to as ‘Optimal Design’.

A lot of physical laws and governing equations can be formulated by means of partial differential equations (PDE) (see [8, 32]). Especially for such kind of equations the finite element method (FEM) is a basic and powerful tool to solve them in an accurate way. The mathematical literature in this field is quite large, and vicarious for the people who work on this topic we refer to [6]. A more engineering literature on FEM is given in [20]. Another method for solving such kind of problems is the so called boundary element method (BEM). We will not go into detail on this topic and refer to [34].

The major problem of FEM is the solution of the arising finite element equations. For reasons of simplicity we restrict ourself on linear ones. If the underlying PDE is elliptic and of second or forth order, multigrid methods (MGM) have proven to be an optimal solver, see [3, 12, 13, 14, 15, 16, 21, 26]. MGMs are known to be asymptotically optimal but they have a great drawback: They require a hierarchical grid structure. All older FEM codes do not work with such kind of underlying grid structure, ergo the MGM could not be implemented in such FEM code as solver. To overcome this problem we can use algebraic multigrid (AMG). This kind of solver preserves the most advantages of MGM and needs no hierarchical grid. Actually we just need the matrix and the right hand side. This makes AMG very attractive as black-box solver with the features of MG-cycling processes.

The first serious approach to AMG was made in 1982 and an improved version can be found in [5]. The algorithm described in there behaves very robustly against parameter jumps, geometry of the domain, etc.; for these problem settings a lot of numerical results can be found in [11, 28, 29]. Additionally the methods described in [7, 35] are also called ‘algebraic’, due to the matrix dependent prolongation and restriction. Anyway they need some hierarchical grid structure. Other ‘pure’ AMG methods with quite interesting approaches can be found in [9, 21, 27, 33]. Again these people are just vicarious for the lot working on this topic.

In the following thesis, we will present AMG methods to get a robust solver for elliptic boundary value problems. Especially we will discuss the method introduced in [5] denoted by RSAMG (Ruge-Stüben algebraic multigrid). For further literature we refer to [28, 29] for a theoretical and practical background; and to [4] for a purely theoretical background.

We will use the solver in three different ways: preconditioner for the conjugate gradient method (CG), iterative solver for a full multigrid method, and black box solver. Particularly we will use RSAMG as preconditioner for CG if we are faced with micro scales. Furthermore we will generalize the algorithm, which is theoretically only for M-matrices, for any symmetric positive definite matrix arising from an FE discretization.

The thesis is organized as follows:

- In Section 2 we consider the mathematical modeling, which results in a very well known scalar PDE, i.e. the potential equation. Accordingly, we give two mathematical methods to reduce the vector valued Maxwell equations to this scalar PDE in 2D. In a straight forward manner a weak formulation is presented in 2D; existence and uniqueness of the solution



are guaranteed.

- Some preparatory work for the FEM is done in Section 3. Here we discuss some possible discretization for 2D problems. Especially we talk about micro scales which may give rise to problems.
- Constructively to Section 3 we present in Section 4 some finite element analysis. There the resulting linear system with its properties, error estimates and a convergence result for the FE solution is included. In this section we also present error estimates for long thin triangles and quadrilaterals.
- The heart of this thesis is Section 5, where a robust AMG solver is presented, i.e. the AMG-method of Ruge and Stüben. We show in that section that such a solver could be used in 3 different forms, i.e. as preconditioner, as iterative solver for FMG, and as black box solver. If we interpret AMG as preconditioner we present a modification to apply RSAMG on all positive definite matrices.
- In Section 6 we demonstrate numerical examples which show the numerical robustness of RSAMG. Additionally coarse grid graphs are shown.
- Last but not least we present in section 7 a summary of this thesis, make some conclusions and pose remarks on a further development.

## 2 Problem description, formulation and preliminaries

In this section we present a very brief discussion of the problem motivation, description and derivation. Therefore a practical problem arising in magnetic field computation is taken. We are concerned with the Maxwell equations and reduce them (for our special problem) to a scalar partial differential equation, i.e. the potential equation (see [22, 20]).

This modeling can be done, at least in 2D areas, with the so called scalar potential or the vector potential. Both ways provide a scalar PDE. The last step in this section is a weak formulation for 2D problems and some mathematical analysis.

### 2.1 A practical problem and its motivation

Let us assume a homogeneous and constant magnetic field (e.g. the magnetic field of the earth in a certain area, assumed to be constant) with induction  $\vec{B}_g = (B_x, B_y)^T = \text{const.}$  A ferromagnetic object with constant permeability  $\mu_F$  is placed in this magnetic field. Moreover we consider only hole objects. (i.e. a so called ‘Faraday’s cage’, see Figure 1). Some typical values for the given parameters are:  $B_x = 21\mu T$ ,  $B_y = 45\mu T$ ,  $\mu_F = 1000$  and  $\mu_{air} = 1$ .

Some points of interest are:

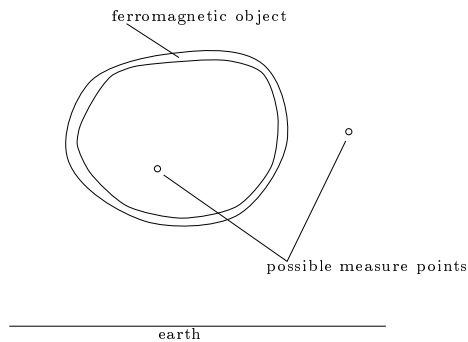


Figure 1: A hole body in the magnetic field of the earth

- (Q1) How ‘big’ is the difference between the undisturbed and the disturbed case in a given measuring point?
- (Q2) How ‘big’ is the magnetic field in the hole of the object?
- (Q3) Can we observe a dependence of the magnetic field on the thickness of the material?
- (Q4) Does the disturbed magnetic field depend on the geometry of the object?
- (Q5) How do material properties influence the magnetic field?

Of course, there could be a lot more questions. Anyway, to answer these questions we have to solve the Maxwell equations, or more exactly a very gentle and special form of them. Then we are able to answer the questions above by some post processing.

Before we start with formulating and modeling we would like to give a brief motivation for such kind of problems. The reason therefore is not only to simulate a ‘Faraday’s cage’. The motivation is stated more in the geometry of the problem. As we will see later, some typical sizes of the problem setting are (see Figure 3): thickness  $t \sim 1 - 4mm$ , diameter of the object  $r \sim 0.5 - 5m$ , and the diameter of the area  $R \sim 2 - 20m$ . Summarized we have a very big area with some micro scales inside and additionally a big parameter jump (from air to ferromagnetic). Both reasons give rise to expect some numerical problems. Therefore we need a very sophisticated solution strategy to solve the problem in finite time and with an appropriate accuracy.

Thus the main topics we would like to discuss in this thesis are

- an appropriate mathematical model and
- a ‘good’ numerical solution strategy.

## 2.2 The physical problem

The basic physical theory are the Maxwell equations for a stationary object (see [22]).

$$\text{rot } \vec{H} = \vec{J} + \frac{\partial \vec{D}}{\partial t} \quad (1)$$

$$\text{rot } \vec{E} = -\frac{\partial \vec{B}}{\partial t} \quad (2)$$

$$\text{div } \vec{B} = 0 \quad (3)$$

$$\text{div } \vec{D} = \rho \quad (4)$$

$$\vec{B} = \underline{\mu} \cdot \vec{H} + \mu_0 \cdot \vec{M}_e \quad (5)$$

$$\vec{D} = \underline{\varepsilon} \cdot \vec{E} + \vec{P}_e \quad (6)$$

$$\vec{J} = \underline{\kappa} \cdot \vec{E} + \vec{J}_e \quad (7)$$

Here  $\vec{M}_e$  is the magnetism of the material,  $\vec{P}_e$  is the electric polarization, and  $\vec{J}_e$  is the impressed current density. These functions could be used as a source in a field problem. The others are:  $\vec{J}$  the current,  $\vec{H}$  the magnetic field,  $\vec{B}$  the magnetic induction,  $\vec{D}$  the electric density and  $\rho$  the charge carrier density.  $\underline{\mu}$ ,  $\underline{\kappa}$  and  $\underline{\varepsilon}$  are time dependent, nonlinear tensors of rank two.

Fortunately, we do not have to solve the whole set of equations, because of the following assumptions.

- we consider a stationary problem, i.e.  $-\frac{\partial \vec{B}}{\partial t} = 0$  and  $\frac{\partial \vec{D}}{\partial t} = 0$
- there is no current, i.e.  $\vec{J} = 0$
- no exciters exist, i.e.  $\vec{M}_e = 0$ ,  $\vec{P}_e = 0$  and  $\vec{J}_e = 0$

Additionally we can assume

- the charge carrier density vanishes, i.e.  $\rho = 0$
- $\mu$  is piecewise constant and bounded from zero, i.e.  $\mu \in L^\infty(\Omega)$
- $\Omega \subseteq \mathbb{R}^2$  or  $\Omega \subseteq \mathbb{R}^3$  is a single connected, bounded domain
- the body  $\Omega_b$  is completely inside the domain  $\Omega$ , i.e.  $\Omega_b \subseteq \Omega$  (see Figure 3)
- there is no residual magnetism in the material

**Remark 2.1.** *Normally the permeability for ferromagnetic materials is a function of the magnetic induction, i.e.  $\mu = \mu(|\vec{B}|)$ . Thus in general the equations are nonlinear. However if the magnetic induction is not too big the assumption of a piecewise, positive function  $\mu$  is a quite good approximation for the real one, see [22, 20].*

Now, in the same way as above we achieve

$$\operatorname{rot} \vec{H} = 0 \quad (8)$$

$$\operatorname{div} \vec{B} = 0 \quad (9)$$

$$\vec{B} = \mu \cdot \vec{H} \quad (10)$$

This set of equations is just slightly different if there is a current, i.e.  $\vec{J} \neq 0$ . With the same assumptions and calculations we get

$$\operatorname{rot} \vec{H} = \vec{J} \quad (11)$$

$$\operatorname{div} \vec{B} = 0 \quad (12)$$

$$\vec{B} = \mu \cdot \vec{H} \quad (13)$$

**Remark 2.2.** *First it is clear that the derived two last systems describe our problem very well in 2D and in 3D. As we will see later, we can always reduce the first system to a scalar PDE (with the scalar potential). Not so the second one. In that case we are only able to reduce it to a scalar PDE in 2D (with the vector potential).*

What is left are boundary and interface conditions (for a more detailed explanation see [22]). Let us assume two materials with the interface  $\Gamma$ , see Figure 2. The interface condition for the magnetic induction and the magnetic density is

$$\vec{B}_1 \cdot \vec{n} = \vec{B}_2 \cdot \vec{n} \quad (14)$$

$$\vec{H}_1 \times \vec{n} = \vec{H}_2 \times \vec{n} \quad (15)$$

respectively, where  $\vec{n}$  denotes the normal unit outward vector. Because the domain  $\Omega$  is finite, we also need boundary conditions for  $\partial\Omega$ .

**Remark 2.3.** *From a mathematical point of view, we could choose at this point a ‘Sommerfeld’s radiation condition’ and couple a BEM with FEM. But for reasons of simplicity we assume a finite domain and thus boundary conditions.*

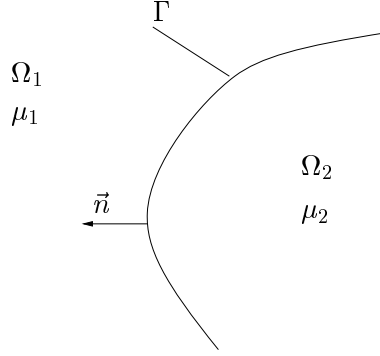


Figure 2: Interface between two different materials

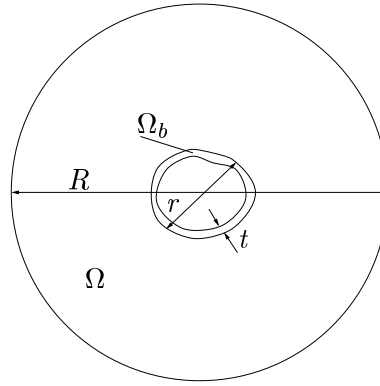


Figure 3: Principle structure of the shielding problem in 2D

It is clear that assuming the domain to be finite would lead to error. But if putting the object  $\Omega_b$  in the domain  $\Omega$  in such a way that the boundary  $\partial\Omega$  is ‘far’ away from the boundary  $\partial\Omega_b$  we can assume the disturbance caused by  $\Omega_b$  to be small (see Figure 3). Thus the undisturbed (i.e.  $\vec{B}_g$ ) should be very close to the disturbed one  $\vec{B}$  at the boundary  $\Gamma$ . (the same is true for the magnetic density). In this way we get the boundary conditions

$$\vec{B} \cdot \vec{n} = \vec{B}_g \cdot \vec{n} \quad (16)$$

$$\vec{H} \times \vec{n} = \vec{H}_g \times \vec{n} \quad (17)$$

for the magnetic induction and the magnetic density, respectively. Again  $\vec{n}$  denotes the unit outward vector.

At this stage, we have an appropriate model of our shielding problem. The next thing we are going to do is to reduce the system of PDEs to a scalar equation.

## 2.3 The mathematical problem

### 2.3.1 The scalar potential

One of the simplest mathematical tricks to reduce the system is to introduce a scalar potential. Because of  $\text{rot } \vec{H} = 0$  we are able to introduce such one for

the magnetic density  $\vec{H}$

$$u : \Omega \rightarrow \mathbb{R} \quad (18)$$

with

$$\vec{H} = -\text{grad } u \quad (19)$$

and consequently

$$\vec{B} = -\mu \cdot \text{grad } u \quad (20)$$

Obviously the existence of such a scalar potential is ensured, at least in a weak sense. A straight forward calculation shows, that we get for  $u(x)$  the scalar equation

$$-\text{div}(\mu(x) \text{grad } u(x)) = 0 \quad \forall x \in \Omega \quad (21)$$

Using (16) for the boundary conditions leads to

$$-\mu(x) \frac{\partial u(x)}{\partial \vec{n}} = \vec{B}_g \cdot \vec{n} \quad \forall x \in \partial\Omega \quad (22)$$

after inserting (16) in (20). With the same arguments we get for the interface condition with (14)

$$-\mu_1(x) \frac{\partial u_1(x)}{\partial \vec{n}} = -\mu_2(x) \frac{\partial u_2(x)}{\partial \vec{n}} \quad \forall x \in \partial\Omega_b \quad (23)$$

According to the boundary conditions, an other approach can also be used. As we can see, a solution for the undisturbed case in 2D is

$$u(x) = -\mu_1(B_x x + B_y y) \quad \forall x \in \Omega \quad (24)$$

Because the solution of the disturbed case should be approximately the solution of the undisturbed case on the boundary  $\partial\Omega$ , the Dirichlet boundary condition, which is attained by (22) and (24)

$$u(x) = -\mu_1(B_x x + B_y y) \quad \forall x \in \partial\Omega \quad (25)$$

can be used. Collecting the results from above we are able to write down the classical formulation of our problem

$$-\text{div}(\mu(x) \text{grad } u(x)) = 0 \quad \forall x \in \Omega_i \quad (26)$$

$$u(x) = -\mu_1(x)(B_x x + B_y y) \quad \forall x \in \partial\Omega \quad (27)$$

$$-\mu_1(x) \frac{\partial u_1(x)}{\partial \vec{n}} = -\mu_2(x) \frac{\partial u_2(x)}{\partial \vec{n}} \quad \forall x \in \partial\Omega_b \quad (28)$$

where  $i \in \{1, 2\}$ .

### 2.3.2 The vector potential

A more general method to reduce (8), (9), (10) or (11), (12), (13) to a scalar equation is to introduce a vector potential. As mentioned above, this can only be done in 2D for a general setting, to achieve a scalar PDE.

Because  $\vec{B}$  is divergence-free and the preliminaries on our domain  $\Omega$  hold we can introduce a vector potential for  $\vec{B}$  (see [10]), i.e.

$$\vec{B} = \text{rot } \vec{A} \quad \text{with} \quad \text{div } \vec{A} = 0 \quad (29)$$

Obviously (9) is always satisfied; thus we get one vector valued equation for the magnetic density

$$\text{rot}\left(\frac{1}{\mu(x)} \text{rot } \vec{A}\right) = 0 \quad (30)$$

An appropriate boundary condition is again (16). Inserting (16) in (29) we obtain for the boundary condition

$$\text{rot } \vec{A} \cdot \vec{n} = \vec{B}_g \cdot \vec{n} \quad (31)$$

and for the interface condition we get similar with (14)

$$\frac{1}{\mu_1} \text{rot } \vec{A}_1 \times \vec{n} = \frac{1}{\mu_2} \text{rot } \vec{A}_2 \times \vec{n} \quad (32)$$

Summarizing, we receive a vector valued equation for  $\vec{A}$  with boundary and interface conditions

$$\text{rot}\left(\frac{1}{\mu(x)} \text{rot } \vec{A}\right) = 0 \quad \forall x \in \Omega_i \quad (33)$$

$$\text{rot } \vec{A} \cdot \vec{n} = \vec{B}_g \cdot \vec{n} \quad \forall x \in \partial\Omega \quad (34)$$

$$\frac{1}{\mu_1(x)} \text{rot } \vec{A}_1 \times \vec{n} = \frac{1}{\mu_2(x)} \text{rot } \vec{A}_2 \times \vec{n} \quad \forall x \in \partial\Omega_b \quad (35)$$

where  $i \in \{1, 2\}$ .

If the 2D case is considered, the above system can be reduced further. For this we assume  $\Omega \subseteq \mathbb{R}^2$ ,

$$\vec{A} = (0, 0, A_3)^T \quad (36)$$

and define  $u(x) := A_3$ . With a little calculation we acquire the scalar PDE

$$-\text{div}\left(\frac{1}{\mu(x)} \text{grad } u(x)\right) = 0 \quad (37)$$

and for the interface and boundary conditions we attain

$$\left\langle \frac{1}{\mu_1(x)} \text{grad } u_1, \vec{n} \right\rangle = \left\langle \frac{1}{\mu_2(x)} \text{grad } u_2, \vec{n} \right\rangle \quad (38)$$

$$-\frac{\partial u}{\partial x} n_1 + \frac{\partial u}{\partial y} n_2 = B_1 n_1 + B_2 n_2 \quad (39)$$

respectively, with  $\vec{B}_g = (B_1, B_2, 0)^T$  and  $\vec{n} = (n_1, n_2, 0)^T$ .

With the same arguments as for the scalar potential a Dirichlet boundary condition can be attained.

$$u(x) = -\frac{1}{\mu_1(x)}(B_y x - B_x y) \quad \forall x \in \partial\Omega \quad (40)$$

Further a scalar PDE of second order with boundary and interface conditions is achieved. Let us collect the results and rewrite them.

$$-\operatorname{div}\left(\frac{1}{\mu(x)} \operatorname{grad} u(x)\right) = 0 \quad \forall x \in \Omega_i \quad (41)$$

$$\left\langle \frac{1}{\mu_1(x)} \operatorname{grad} u_1, \vec{n} \right\rangle = \left\langle \frac{1}{\mu_2(x)} \operatorname{grad} u_2, \vec{n} \right\rangle \quad \forall x \in \partial\Omega \quad (42)$$

$$u(x) = -\frac{1}{\mu_1(x)}(B_y x - B_x y) \quad \forall x \in \partial\Omega_b \quad (43)$$

**Remark 2.4.** *Another quite interesting problem is the following:*

*Let us assume that the time dependency is not neglectable, then we get with the help of the vector potential the equation*

$$-\operatorname{div}\left(\frac{1}{\mu} \operatorname{grad} u(x, t)\right) = -\kappa \frac{\partial u(x, t)}{\partial t} - \epsilon \frac{\partial u(x, t)^2}{\partial t^2} \quad (44)$$

where  $x$  is the space and  $t$  the time variable. Suppose  $u(x, t)$  to be a time-harmonic, i.e. it is of the form

$$u(x, t) = e^{i\omega t} \cdot U(x) \quad (45)$$

and let us consider the damping term  $-\kappa \frac{\partial u(x, t)}{\partial t}$  to be neglectably small. Then we attain for the space an elliptic problem, i.e. the Helmholtz equation.

$$-\operatorname{div}\left(\frac{1}{\mu(x)} \operatorname{grad} U(x)\right) + \omega^2 U(x) = 0 \quad (46)$$

## 2.4 The resulting variational form

### 2.4.1 A weak formulation in 2D

Let us consider (26), (27), (28) and a geometrical setting like in Figure 3 as basis for our weak formulation. Without going into detail for weak formulations we give just the result. For the theory on Sobolev spaces see [1]; for weak formulations see [6, 25].

An appropriate function space for elliptic PDEs of second order is the  $H_0^1(\Omega)$ . Further we set  $V := H_0^1(\Omega)$  and

$$V_g := \{v \in V : v - g|_{\partial\Omega} = 0\} \quad (47)$$

where  $g$  is the given Dirichlet boundary condition on  $\partial\Omega$ . In a straight forward manner the following bilinear and linear form is obtained after homogenization.

$$a(u, v) := \int_{\Omega} \mu(x) \operatorname{grad} u(x) \operatorname{grad} v(x) dx \quad \forall v \in V \quad (48)$$

$$\langle F, v \rangle := -a(g, v) \quad \forall v \in V \quad (49)$$



respectively, which can be written as abstract variational problem

$$\text{find } u \in V : a(u, v) = \langle F, v \rangle \quad \forall v \in V \quad (50)$$

### 2.4.2 Existence and uniqueness

The basic tool for a satisfactory existence and uniqueness theory is the Theorem of Lax-Milgram. First the general result is shown and afterwards an existence and uniqueness result is attained for our counter example.

**Theorem 2.5.** (*Lax-Milgram*) *Let  $F \in V^*$  and  $a(.,.) : V \times V \rightarrow \mathbb{R}$  be a bilinear form which fulfills*

- *V-elliptic* :  $\exists \zeta_1 : \zeta_1 \|v\|^2 \leq a(v, v) \quad \forall v \in V$
- *V-continous* :  $\exists \zeta_2 : |a(u, v)| \leq \mu_2 \|u\| \|v\| \quad \forall u, v \in V$

*then there exists a unique solution of*

$$\text{find } u \in V : a(u, v) = \langle F, v \rangle \quad \forall v \in V \quad (51)$$

*Proof.* can be found in [24] □

**Corollary 2.6.** *The equation (50) has a unique solution in  $V$ .*

*Proof.*  $F \in V^*$  and  $a(u, v)$  is a bilinear form are obviously fulfilled. Next the V-ellipticity is shown.

$$\begin{aligned} a(u, u) &= \int_{\Omega} \mu(x) \text{grad } u(x) \text{grad } u(x) dx \\ &\geq \min_{i \in \{1,2\}} \{\mu_i\} \int_{\Omega} \text{grad } u(x) \text{grad } u(x) dx \\ &= \min_{i \in \{1,2\}} \{\mu_i\} (0.5 \int_{\Omega} \text{grad } u(x) \text{grad } u(x) dx + 0.5 \int_{\Omega} \text{grad } u(x) \text{grad } u(x) dx) \\ &\geq 0.5 \min_{i \in \{1,2\}} \{\mu_i\} \min \{1, c_F\} \int_{\Omega} \text{grad } u(x)^2 + |u|^2 dx \\ &= \zeta_1 \|u\|^2 \end{aligned} \quad (52)$$

with  $\zeta_1 = 0.5 \cdot \min_{i \in \{1,2\}} \{\mu_i\} \cdot \min \{1, c_F\}$  and  $c_F$  the constant of the Friedrichs inequality. To finish the proof we need the V-continuity.

$$\begin{aligned} |a(u, v)| &= \int_{\Omega} \mu(x) \text{grad } u(x) \text{grad } u(x) dx \\ &\leq \max_{i \in \{1,2\}} \{\mu_i\} \int_{\Omega} \text{grad } u(x) \text{grad } u(x) dx \\ &\leq \max_{i \in \{1,2\}} \{\mu_i\} \sqrt{|\text{grad } u|^2} \sqrt{|\text{grad } v|^2} \\ &\leq \zeta_2 \|u\| \|v\| \end{aligned} \quad (53)$$

□

**Remark 2.7.** *Because the magnetic induction can be represented by*

$$\vec{B} = -\text{grad } u(x) \quad \forall x \in \Omega \quad (54)$$

*the solution of the magnetic induction is unique.*

### 2.4.3 Weak formulation for the vector-potential

A FEM discretization is much more delicate if we consider the reduced Maxwell equations (11), (12), (13) in 3D. Why is it difficult?

Normally, we can use for our FE-discretization Lagrange elements because of their well known features if the continuous problem is formulated in  $H^1(\Omega)$ . For a general magnetic field problem in 3D an appropriate function space is  $H_0(\text{rot}) \cap H(\text{div})$  for the weak formulation. The difficulty is that  $H^1(\Omega) \subset H_0(\text{rot}) \cap H(\text{div})$  if the domain  $\Omega$  is not convex. A consequence therefore is that we are not allowed to discretize  $H_0(\text{rot}) \cap H(\text{div})$  with Lagrange elements, because this could lead to a completely wrong solution. Some theory on this topic can be found in [18, 19, 23].

### 3 Discretization in 2D

A discretization of a domain  $\Omega$  has to be carried out very carefully if we want to use it building an FE-space. For this reason we give a brief recapitulation of what we mean by a regular triangulation first. Consequently we also provide a definition for a regular quadrilateral discretization. Secondly we discuss the problem of micro scales and accordingly we present some possible methods of discretizing them. In the end of this section we suggest refinement strategies for long thin triangles and quadrilaterals. Especially this last topic is of great interest for an adaptive refinement strategy.

#### 3.1 General remarks

Before we are able to build up a finite element space we have to discretize our domain  $\Omega$ . This can be done in many ways but we will just concentrate on the discretization with triangles and quadrilaterals.

As usual we call a family  $(\tau_h)_{h \in \Theta}$ , with  $\tau_h := \{\delta_r : r \in \mathbb{R}_h\}$ , of triangles and quadrilaterals a regular discretization of  $\Omega$  if for all  $t, \bar{t} \in \tau_h$  with  $t \neq \bar{t}$  the intersection of them is either a common vertex, a common edge, or empty (for examples see Figure 4). Furthermore we have to assume some geometrical features on the elements of the discretization. Therefore we denote by

$$h_{max}^i = \max_j h_j^i \quad (55)$$

$$h_{min}^i = \min_j h_j^i \quad (56)$$

where  $h^i$  denotes the set of edges in an element  $i$ . The subscript  $j$  is either  $j = 1, 2, 3$  for triangles or  $j = 1, 2, 3, 4$  for quadrilaterals. Moreover we abbreviate

$$h_{max} = \max_{i \in \mathbb{R}_h} h_{max}^i \quad (57)$$

$$h_{min} = \min_{i \in \mathbb{R}_h} h_{min}^i \quad (58)$$

the longest and the shortest edge in the discretization, respectively.

For triangles we assume the smallest angle of all triangles in  $\tau_h$  to be uniformly bounded away from zero and

$$\frac{h_{max}^i}{h_{min}^i} \leq c_t < \infty \quad (59)$$

holds for every triangle  $i$ , with a constant  $c_t > 0$ , independent of  $i$ . In the same manner we assume for a quadrilateral

$$|\cos \theta_k| \leq \sigma < 1 \quad (60)$$

with  $\theta_k$ ,  $k = 1, 2, 3, 4$ , the inner angles of the quadrilateral and

$$\frac{h_{max}^i}{h_{min}^i} \leq c_q < \infty \quad (61)$$

holds for every quadrilateral  $i$ , with a constant  $c_q > 0$ , independent of  $i$ . As we can easily see there are also long thin triangles and rectangles included if we allow the constants  $c_t, c_q$  to be relatively large. Finally we abbreviate the set of inner nodes, all nodes and edges by  $\omega_h, \bar{\omega}_h$  and  $E_h$ , respectively.  $\Omega_h$  denotes the union of all triangles in  $\tau_h$ .

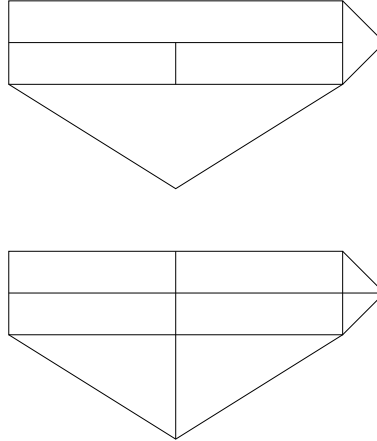


Figure 4: A non conform and a conform discretization

### 3.2 The problem of micro scales

Before we start to discuss an appropriate discretization method we need to know what we mean by a micro structure in 2D.

Let us assume an area  $\Omega$  which is subdivided by subareas  $\Omega_i$  in such a way, that each  $\Omega_i$  has the same constant permeability  $\mu_i$  ( if there is a hole in  $\Omega$  we assume the density function there to be zero). This is necessary because a discretization should not cross an interface line.

Next heuristics are given to decide whether a subdomain is a micro structure or not; for examples of typical micro structures in  $\mathbb{R}^2$  see Figure 5.

- A subregion  $\Omega_i$  is called a micro structure of type 1 if the diameter of the area is small compared to the diameter of  $\Omega$ .
- A subregion  $\Omega_i$  is called a micro structure of type 2 if the diameter of the area is big compared to the smallest measurement, e.g. the thickness.

It is quite clear that micro structures of type 1 need no further attention because we hardly have any other possibility for discretizing them than with normal triangles or quadrilaterals, i.e.  $c_t \sim 2$  or  $c_q \sim 2$ .

A micro structure of type 2 has to be concerned more carefully. This can be seen in the following. Let us consider an area like in Figure 6. If we discretize micro structure in a naive way we obtain a discretization which has one order more unknowns as if we use long thin triangles or quadrilaterals. This may be crucial in the context of FEM because the arising linear system becomes considerably larger in the naive way than in the more sophisticated way using long thin elements.

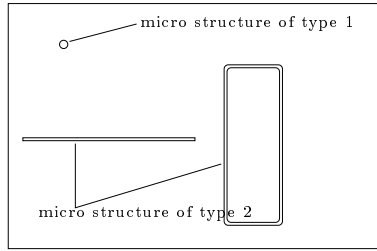


Figure 5: Typical micro structures in 2D

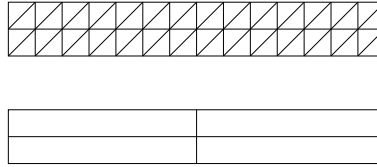


Figure 6: Two discretizations of a micro structure of type 2

### 3.3 Refinement strategy

If we would like to use an adaptive technique or simple would like to refine the whole discretization one time it is necessary to have an appropriate refinement strategy for the given discretization. For triangles there is hardly another useful method than quartering in four congruent smaller triangles (see Figure 7). It

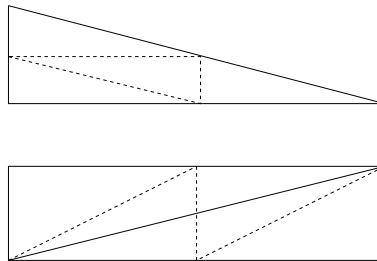


Figure 7: Refinement strategies for long thin triangles

can be easily seen that the smallest angle in the triangle can not be improved with some technique. Unfortunately a very small angle implies a large condition number; ergo long thin triangles are not a very useful tool for micro structures of type 2.

An improved method is as follows. Suppose a quadrilateral, which is divided into 2 triangles, so that the biggest angle in the quadrilateral is divided into two smaller ones (see Figure 7). A refinement strategy is to see the two triangles as the original quadrilateral and first half the quadrilateral and then to make two triangles. This can be done till  $h_{min}^i \sim h_{max}^i$ . In this way the smallest angle is improved without shrinking the smallest edge of the discretization.

According to the last technique a refinement strategy is obtained for a quadri-

lateral by a simple cut of the two longest sides (see Figure 8) till the smallest and the longest side are of the same order, i.e.  $h_{min}^i \sim h_{max}^i$ . If  $h_{min}^i \sim h_{max}^i$  then two triangles are made of the quadrilateral.

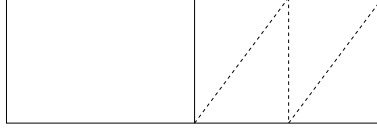


Figure 8: Refinement strategy for long thin rectangles

**Remark 3.1.** *It still could be useful to quarter a quadrilateral, i.e. for an anisotropic operator. Anyway, for our demands the strategy above seems to be the best one, because in this way the number of unknowns is kept reasonably small.*

## 4 Finite element analysis

Constructive to Section 3 we build up an FE-space in this section. The space will be spanned by linear and bilinear form-functions. For this kind of setting we give the most famous properties of the resulting linear system. Especially we show estimates for the condition-number in the case of long thin elements. Next we present a convergence result in  $H^1$ . For that, estimates for the approximation error for long thin elements are shown, if the quotient of the inner and the outer radius of the element is uniformly bounded from above by a constant.

### 4.1 The resulting discrete problem - general remarks

#### 4.1.1 The finite element space

Before we are able to make a convergence theory we have to construct an appropriate finite dimensional function space. For this we use the discretization introduced in the last section and define linear form-functions on triangles and bilinear form-functions on quadrilaterals. Thus we have introduced nothing else than the nodal basis, which is characterized by

$$\phi_i(x) = \begin{cases} 1 & \text{if } x = p_i \\ 0 & \text{else} \end{cases} \quad (62)$$

where  $\phi_i(\cdot)$  is the  $i^{\text{th}}$  base function and  $p_i \in \omega_h$  is a vertex. Obviously the FE-space is defined by

$$V_h = \{\phi_i : p_i \in \bar{\omega}_h\} \quad (63)$$

$$V_{0h} = \{\phi_i : p_i \in \omega_h\} \quad (64)$$

For more details how to build an FE-space we refer to [6, 25, 20].

Now the following very important result holds.

**Theorem 4.1.** *Let  $V_h \subseteq H^1(\Omega_h)$  be the finite element space introduced in (63) then we get*

$$V_h \cap H^1(\Omega_h) \subseteq C^0(\Omega_h) \quad (65)$$

*Proof.* can be found in [6] □

By using the FE-space defined in (63), we get the following equation for the finite element approximation.

$$a(u_h, v_h) = \langle F, v_h \rangle \quad \forall v_h \in V_h \quad (66)$$

where the bilinear form and the linear form are defined as in Section 2 by the equations (48) and (49), respectively, and  $u_h$  is the finite element approximation. In this way we obtain a linear system

$$A_h \underline{u}_h = \underline{f}_h \quad (67)$$

where the underlined variables denote the coefficient vector of a function in  $V_h$ . Next we want to list some properties for our resulting linear system.

- 1) The matrix  $A_h$  is an M-matrix, except if we use triangles where an angle is greater than  $\pi/2$  or long thin quadrilaterals.
- 2) The resulting matrix is symmetric and positive definite if we use the vector potential with Dirichlet boundary conditions.
- 3) The matrix is large, sparse and poor conditioned.

Particularly the first property is of special interest for us, because the algorithm presented in Section 5 to solve (67) is just well suited for M-matrices. Consequently we have to adapt the algorithm if the M-matrix property is lost. Obviously the resulting linear system (67) strongly depends on the underlying grid.

#### 4.1.2 The condition number

Let us consider the condition number of such a linear system more precisely. The following theorem shows the behavior of the condition number if we do not use long thin elements.

**Theorem 4.2.** *Let us assume the standard preliminaries (see Section 2) on the variational problem and suppose the triangulation to be regular, i.e.  $c_t, c_q \sim 5$ , thus  $h_{min} \sim h_{max} \sim h$ .*

*Then we get the following sharp estimate for the eigenvalues. There exist constants  $\bar{c}_E, \underline{c}_E > 0$  independent of  $h$  such that*

$$\underline{c}_E h^2 \leq \lambda(K_h) \leq \bar{c}_E h^0 \quad (68)$$

*and straightforward we get for the condition number*

$$\kappa(K_h) \leq \frac{\bar{c}_E}{\underline{c}_E} h^{-2} \quad (69)$$

*where  $h$  is the global mesh size.*

*Proof.* can be found in [25] □

Supposing a discretization with long thin triangles and rectangles with the features introduced in point 3.1 we attain

**Theorem 4.3.** *Let us assume the standard preliminaries on the variational problem and suppose a regular discretization with  $c_t, c_q \gg 1$ .*

*Then we get the following estimate for the eigenvalues. There exist constants  $\bar{c}_E, \underline{c}_E > 0$  independent of  $h_{min}, h_{max}$  such that*

$$\underline{c}_E h_{min}^2 \leq \lambda(K_h) \leq \bar{c}_E h_{max}^0 \quad (70)$$

*and straightforward we get for the condition number*

$$\kappa(K_h) \leq \frac{\bar{c}_E}{\underline{c}_E} c_h h_{min}^{-2} \quad (71)$$

*with  $c_h = \frac{h_{max}}{h_{min}}$*



*Proof.* It is very similar to the last theorem.  $\square$

To illustrate this result we assume two discretizations. The first one with a global discretization parameter  $h = 10^{-2}$  and the second one with long thin elements where  $h_{min} = 10^{-2}$  and  $h_{max} = 1$ . The constant arising for long thin elements additionally is  $c_h = 1000$ . Concluding a condition number of  $10^4$  and  $10^7$  has to be expected for the first and the second example, respectively.

## 4.2 The approximation error and a convergence result

### 4.2.1 The approximation error

Without going into detail, we give an error estimate for triangles and quadrilaterals in  $H^1(\Omega)$  first. For error estimates on triangles we refer to [6, 25]. For the approximation error estimates of quadrilaterals we follow [36]. For more sophisticated estimates for anisotropic elements we refer to [2].

**Theorem 4.4.** *Let us assume*

- (1)  $\Omega \subset \mathbb{R}^m$  be a bounded domain with  $\partial\Omega \in C^{0,1}$  with a regular triangulation (long thin triangles are allowed)
- (2)  $F(\Delta) \supset P_k(\Delta)$
- (3)  $u \in V_g$  and  $u \in W_2^{k+1}(\delta_r) \forall r \in \mathbb{R}_h \forall h \in \Theta$

Above  $F(\Delta)$  and  $P_k(\Delta)$  are abbreviations for the space of linear form functions and the space of polynomials with degree less equal  $k$ , respectively.

Then there exists a constant  $c > 0$ , so that

$$\inf_{v_h \in V_h} |u - v_h|_{1,\Omega} \leq c_1 h_{max}^k \left( \sum_{r \in \mathbb{R}_h} |u|_{k+1}^2 \right)^{1/2} \quad (72)$$

*Proof.* A proof for a regular triangulation  $c_t, c_q \sim 5$  can be found in [25]. Consequently a proof with long thin elements is very similar if they have the property  $c_t, c_q < \infty$ .  $\square$

Of our special interest is the case  $k = 1$ , i.e. linear form-functions. Thus the approximation error is  $O(h_{max})$ . For other types of elements, i.e. long thin quadrilaterals, we pose a similar result for bilinear form-functions.

**Theorem 4.5.** *Let  $K$  be a narrow quadrilateral with parallel long sides which satisfies  $h_{min} \leq \frac{1}{12} h_{max}$  (see Figure 9). Let  $u \in W_2^2(K)$ . Then the following estimate holds*

$$\|u - u_h\|_{1,K} \leq c h_{max} |u|_{2,K} \quad (73)$$

*Proof.* can be found in [36]  $\square$

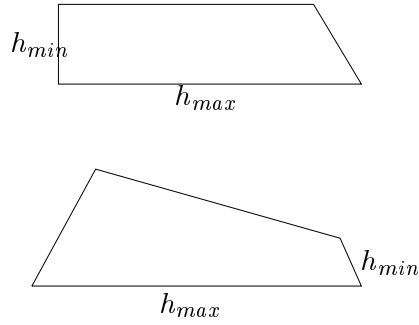


Figure 9: Two possible quadrilateral elements

In [36] they show in straight forward manner a generalization of the last theorem in the sense if no parallel sides exist in a quadrilateral (see Figure 9). Under sufficient conditions the estimate

$$\|u - u_h\|_{1,K} \leq c_2 h_{max} |u|_{2,K} \quad (74)$$

is attained.

**Remark 4.6.** *In literature a lot of other error estimates are known, e.g. in  $L_\infty$  or  $W_\infty^1$ . We give here no further results on this topic and refer to [6, 25] and references therein.*

#### 4.2.2 A convergence result in $H^1$

With the approximation results for long thin triangles and quadrilaterals above, we are able to present a convergence result for the FE solution in  $H^1$ .

**Theorem 4.7.** *Let us assume*

- (1) *standard preliminaries on the variational form*
- (2)  *$\Omega \subset \mathbb{R}^2$  bounded domain with  $\partial\Omega \in C^{0,1}$  which is regular triangularized in the sense of Point 3.1*
- (3)  *$F(\Delta) \supset P_1(\Delta)$*
- (4)  *$V_{gh} = g_h + V_{0h} \subset V_g$  with  $V_{0h} \subset V_0$  and given  $g_h \in V_g \cap V_h$*
- (5)  *$u \in V_g : a(u, v) = \langle F, v \rangle$  for all  $v \in V_0$  and  $u_h \in V_{gh} : a(u_h, v_h) = \langle F, v_h \rangle$  for all  $v_h \in V_{0h}$*
- (6)  *$u \in V_g$  and  $u \in W_2^2(\delta_r) \forall r \in \mathbb{R}_h \forall h \in \Theta$*

*Then we get*

$$\|u - u_h\|_1 \leq c \cdot h_{max} \cdot \left( \sum_{r \in \mathbb{R}_h} |u|_2^2 \right)^{1/2} \quad (75)$$

*Proof.* With Theorem 4.4 and 4.5 we get immediately

$$\begin{aligned} \|u - u_h\| &\leq \begin{cases} c_1 h_{max} (\sum_{r \in \mathbb{R}_h^t} |u|_2^2)^{1/2} \\ c_2 h_{max} (\sum_{r \in \mathbb{R}_h^q} |u|_2^2)^{1/2} \end{cases} \\ &\leq \max \{c_1, c_2\} h_{max} (\sum_{r \in \mathbb{R}_h} |u|_2^2)^{1/2} \end{aligned} \quad (76)$$

where  $\mathbb{R}_h^t$  and  $\mathbb{R}_h^q$  denote the set of triangles and quadrilaterals, respectively.  $\square$

## 5 Numerical solution

In the previous sections we discussed for our problem an appropriate discretization of our domain and did some analysis of the resulting finite element problem. What is left is a practical solver for the arising linear, sparse and poor conditioned system. Additionally the technique to solve such matrices should be robust against parameter jumps, micro scales and the geometry. Further we have to face another problem: The number of unknowns on the coarsest grid is of order  $10^3 - 10^5$  and it is not reduceable without loss of essential information. In this section we give one way out of these problems and present a robust method, i.e. algebraic multigrid of Ruge and Stüben, abbreviated as RSAMG. In the first subsection we present three areas where RSAMG can be used, i.e. as black box solver, as iterative solver for the full multigrid method, or as preconditioner for the conjugate gradient method. Especially the last technique leads to a solver demanded above. Next we are going to describe and motivate the RSAMG algorithm. This algorithm is understood best and performs very well in the case of M-matrices. Finally we present a method to solve general s.p.d. matrices with the same algorithm as described in this section, if the matrix is received from an FE-discretization.

### 5.1 Motivation and practical usage of RSAMG

As mentioned in the brief introduction of this section there are a lot of demands on the solver. Additionally the solver should be of optimal order. One way to solve a matrix equation is the general and very flexible setting of an AMG method, which is very robust and can still be programmed in  $O(n)$ ; here and later on  $n$  is the dimension of the matrix. Before we discuss the technique of RSAMG three examples are treated how RSAMG can be used and motivated.

#### 5.1.1 RSAMG as a black box solver

Actually RSAMG needs just the matrix and the right hand side. In addition to that it is clear that RSAMG is necessarily less efficient than normal MG. However there are a lot of problems, where normal MG breaks down and RSAMG performs still very well. Let us present some problems where geometrical MG-methods could break down.

- (1) Let us consider a domain  $\Omega$  with some regular triangulation and we would like to solve an elliptic PDE. The resulting system-matrix is s.p.d., but a normal MG-method can not be applied if there is no hierarchical grid underlying the discretization. For RSAMG such kind of problems are no difficulty and this property makes it very attractive as solver in an older FE-code.
- (2) The next kind of problem we will discuss are those where the domain is complex enough, so that any sensible discretization is too fine to serve as coarse grid (the number of unknowns on the coarsest grid is too big). Thus the memory requirement and the solution time of a classical iterative

method are too high. For the same reasons finer grids are not possible. Again RSAMG has no problem with such kind of discretization, because coarsening is purely automatic and does not depend on the underlying geometry.

- (3) This kind of problem is closely related to the second one and appears if no uniform coarsening is available, i.e. we have to use irregular grids. Particularly for our counter examples we take long thin elements; this could be an example for that problem class.
- (4) A very interesting problem in means of FEM is not related to the domain but to the operator. This is for example an anisotropic operator or the convection-diffusion equation. Moreover there could be discontinuous coefficients or singular perturbed equations.
- (5) Last but not least there are purely discrete problems, with no geometrical background. Also in this case RSAMG perform excellent.

The principle setting of such a solver is very easy. In contrast to geometrical MG-methods the smoothing operator is fixed for every problem and usually the simplest one is taken, i.e. Gauß-Seidel or Jacobi smoother. Additionally a sophisticated coarsening strategy is used, which is closely related to the interpolation weights. Finally the coarse grid operator is defined via the Galerkin identity. This strategy can always be applied for s.p.d. M-matrices.

### 5.1.2 Iterative solver for FMG

One very famous solution strategy is the so called full multigrid (FMG) method. It is based on the nested-iteration principle and the most important feature is: the method is of optimal order. In the next theorem we show (under sufficient preliminaries on the MG-operator), that it is enough to solve the coarsest grid not exactly but just with a given accuracy (to get the same accuracy for the solution). This is important if the coarsest grid of an FMG-cycle could not be reduced geometrically and has too many unknowns to solve it directly. The next theorem shows us how to handle this problem practically.

**Theorem 5.1.** *Let  $M_q$  be the multigrid operator of level  $q$ . Further on let  $\underline{u}_{q-1}$  and  $\underline{u}_q$  be the exact solutions of level  $q$  and  $q-1$  for  $A_q \underline{u}_q = \underline{f}_q$  and let  $I_{q-1}^q \underline{u}_{q-1}$  the interpolated solution of level  $q$ . Assuming that the following preliminaries hold*

- $\|M_q\| \leq \eta < 1$
- $\|I_{q-1}^q \underline{u}_{q-1}\| \leq c_1(u) h_q^p$
- $\|I_{q-1}^q\| \leq c_I$
- $\frac{h_{q-1}}{h_q} \leq \alpha$

where all constants are positive and are not depending on  $h_q$ .

Moreover let  $k_q = i$  for all  $q = 2, \dots, l$  with  $i: c_2 \eta^i < 1$ , and  $c_2 = c_I \alpha^p$ .

If there are  $i$  MG-iteration steps applied on every nested level, then we get for  $\underline{v}_q^i$  (the FMG approximation) the following estimate holds

$$\|\underline{u}_q - \underline{v}_q^i\| \leq c_3(\eta) c_1(u) h_q^p \quad (77)$$

with  $c_3(\eta) = \frac{\eta^i}{1 - c_2 \eta^i}$ .

*Proof.* can be found in [17] □

If we are faced with a strong elliptic PDE of second order the arising constants in Theorem 5.1 can be determined more exactly.

**Theorem 5.2.** *Let us suppose the assumptions of Theorem 5.1 and  $p = 1$ , i.e. linear form functions and the estimate*

$$\|u - u_q\| \leq \frac{\zeta_1}{\zeta_2} \cdot c \cdot h_q \cdot \left( \sum_{r \in \mathbb{R}_h} |u|^2 \right)^{1/2} \quad (78)$$

Then we get

$$\|\underline{u}_q - \underline{v}_q\| \leq \frac{\eta^i}{1 - \alpha \eta^i} c_1(u) h_q \quad (79)$$

$$\|u_q - v_q\| \leq \frac{\zeta_2}{\zeta_1} c |u|_2 \left( 1 + 3 \sqrt{\frac{\zeta_2}{\zeta_1} \frac{\eta^i}{1 - \alpha \eta^i}} \right) \quad (80)$$

with  $c_1(u) = 3 \sqrt{\frac{\zeta_2}{\zeta_1}} c |u|_2$  and  $c_I = 1$ ,  $\alpha \sim 2$

*Proof.* can be found in [17] □

It is easy to see how to find a good solution strategy. We are able to use an iterative solver for the coarse grid ( $q = 1$ ), i.e. we are allowed to use RSAMG as coarse grid solver.

For an algorithm see Algorithm 1

```

PROCEDURE CoarseGridSolver( $A_l, \underline{f}_l$ )
BEGIN

  WHILE  $\|\underline{u}_1 - \underline{v}_1^i\| \geq c_3(\eta) c_1(u) h_1^p$  DO
    BEGIN
      RSAMG( $A_l, \underline{f}_l$ )
    END

END

```

Algorithm 1: Strategy for an iterative coarse grid solver

### 5.1.3 RSAMG as a preconditioner

As mentioned above we can use or interpret RSAMG as a preconditioner (for CG). For an algorithm of a preconditioned CG method see [25, 17]. This combination will result in a very robust solution strategy as will be shown numerically in the next section. But before, we present the main results on this technique. For more details we refer to [17].

Before we start we have to introduce some notations and abbreviations.

**Definition 5.3.** *Let  $A, B \in \mathbb{R}_n^n$  be s.p.d., then we call  $B$  spectrally equivalent to  $A$  with spectral constants  $0 < \underline{c}_B, \underline{c}_B < \infty$  if for all  $u \in \mathbb{R}^n$  the following inequalities hold*

$$\underline{c}_B(Bu, u) \leq (Au, u) \leq \bar{c}_B(Bu, u) \quad (81)$$

Abbreviated we will write for (81) simply  $\underline{c}_B B \leq A \leq \bar{c}_B B$  or  $A \sim B$ .

The first theorem provides the spectral constants of a MG preconditioner applied to the system matrix.

**Theorem 5.4.** *Let us assume that*

- $B_l$  s.p.d. with  $\underline{\gamma}_B B_l \leq K_l \leq \bar{\gamma}_B B_l$
- $B_q$  s.p.d. for all  $q = 1, \dots, l-1$
- $S_q^{post} = (S_q^{pre})^{*B_q}$  for all  $q = 2, \dots, l$
- $I_q^{q-1} = (I_{q-1}^q)^T$  for all  $q = 2, \dots, l$

Moreover we assume for the MG-operator the estimate

$$\|M_l\|_{B_l} \leq \eta < 1 \quad (82)$$

where  $\eta$  is a constant independent of  $h_l$ .

Then the MG-preconditioner

$$C_l = B_l(I_l - (M_l)^k)^{-1} \quad (83)$$

is symmetric and positive definite. Moreover we get the spectral inequalities

$$\underline{\gamma}_C C_l \leq A_l \leq \bar{\gamma}_C C_l \quad (84)$$

with

$$\underline{\gamma}_C = \underline{\gamma}_B(1 - \eta^k) \quad (85)$$

$$\bar{\gamma}_C = \begin{cases} \bar{\gamma}_B(1 + \eta^k), & \text{if } k \text{ is odd} \\ \bar{\gamma}_B, & \text{if } k \text{ is even} \end{cases} \quad (86)$$

*Proof.* can be found in [17] □

**Remark 5.5.** In the last and the next theorems  $\eta$  is also called the MG-rate. Moreover the index  $l$  denotes the finest level and  $I_n$  is the abbreviation for the identity-operator in  $\mathbb{R}^n$ . Actually there are two major preliminaries in the last theorem, which should be mentioned

- (i) The condition  $S_q^{post} = (S_q^{pre})^{*B_q}$  for all  $q = 2, \dots, l$  is fulfilled for smoothing operators of interest, i.e. Jacobi, Gauß-Seidel forward/backward.
- (ii) For RSAMG we can not show the estimate (82) for a multilevel cycling process. We just present an estimate for the two grid method. Nevertheless the practical examples seem to confirm (82) also in the multilevel case.

According to the last theorem we present a slightly more general theorem as tool to improve the spectral constants in (85) and (86).

**Theorem 5.6.** Let us assume that

- $B_l$  s.p.d. with  $\underline{\gamma}_B B_l \leq K_l \leq \bar{\gamma}_B B_l$
- $B_q$  s.p.d. for all  $q = 1, \dots, l-1$
- $S_q^{post} = (S_q^{pre})^{*B_q}$  for all  $q = 2, \dots, l$
- $I_q^{q-1} = (I_{q-1}^q)^T$  for all  $q = 2, \dots, l$

and moreover the inequalities

$$-\eta_1 \leq \mu(M_l) \leq \eta_2 \quad (87)$$

hold, where  $\eta_1, \eta_2 \in [0, 1)$ , and  $\mu(M_l)$  is the set of eigenvalues of  $M_l$ . Then the MG-preconditioner

$$C_l = B_l(I_l - (M_l)^k)^{-1} \quad (88)$$

is symmetric and positive definite. Moreover we get the spectral inequalities

$$\underline{\gamma}_C C_l \leq A_l \leq \bar{\gamma}_C C_l \quad (89)$$

with the constants

$$\underline{\gamma}_C = \underline{\gamma}_B (1 - \eta_2^k) \quad (90)$$

$$\bar{\gamma}_C = \begin{cases} \underline{\gamma}_B (1 + \eta_1^k), & \text{if } k \text{ is odd} \\ \bar{\gamma}_B, & \text{if } k \text{ is even} \end{cases} \quad (91)$$

*Proof.* can be found in [17] □

Finally we give a theorem which shows, that  $\bar{\gamma}_C$  could be improved if we determine the a-priori preconditioners  $B_q$  with the Galerkin formula

**Theorem 5.7.** Let us assume that

- $B_l$  s.p.d. with  $\underline{\gamma}_B B_l \leq K_l \leq \bar{\gamma}_B B_l$



- $B_q$  s.p.d. for all  $q = 1, \dots, l-1$
- $S_q^{post} = (S_q^{pre})^{*B_q}$  for all  $q = 2, \dots, l$
- $I_q^{q-1} = (I_{q-1}^q)^T$  for all  $q = 2, \dots, l$
- $I_q^{q-1}$  have full rank for all  $q = 2, \dots, l$
- $I_q^{q-1} B_q I_{q-1}^q \leq B_{q-1}$

Moreover we assume for the MG-operator the estimate

$$\|M_l\|_{B_l} \leq \eta < 1 \quad (92)$$

where  $\eta$  is a constant independent of  $h_l$ .

Then the MG-preconditioner

$$C_l = B_l(I_l - (M_l)^k)^{-1} \quad (93)$$

is symmetric and positive definite. Moreover we get the spectral inequalities

$$\underline{\gamma}_C C_l \leq A_l \leq \bar{\gamma}_C C_l \quad (94)$$

with the spectral equivalence constants

$$\underline{\gamma}_C = \underline{\gamma}_B(1 - \eta^k) \quad (95)$$

$$\bar{\gamma}_C = \bar{\gamma}_B \quad (96)$$

*Proof.* can be found in [17] □

**Remark 5.8.**

(i) In practice we use mostly  $B_l = A_l$ , therefore we could easily compute the spectral constants. However, as the theorem shows, it is admissible to use a spectrally equivalent matrix  $B_l$  to  $A_l$  as an a-priori preconditioner. This fact will be very important for us, because in this way RSAMG can be applied not only to s.p.d. M-matrices.

(ii) The preliminaries on the a-priori preconditioner and on the interpolation operator hold for RSAMG. Unfortunately we can give no proof which ensures (82).

## 5.2 AMG of Ruge and Stüben (RSAMG)

Now we are going to describe the method of Ruge and Stüben. Therefore we will follow mainly [28].

### 5.2.1 Terminology, assumptions and notations

In order to solve a (sparse) linear system

$$A_h \underline{u}_h = \underline{f}_h \quad (97)$$

by means of a multigrid-like cycling process, we have to introduce the major tools for a two grid cycle, i.e.

- coarse grid operator:  $A_H$
- interpolation operator:  $I_H^h$
- restriction operator:  $I_h^H$
- smoothing operator:  $S_h$

where the subscript  $H$  denotes the coarse-grid and  $h$  the finer grid. If we use  $q$  and  $q + 1$  instead of  $h$  and  $H$ , respectively the notation is meant in the context of multigrid cycles. Once the two grid operators are known we can set up a multigrid cycle in the usual way (see [17]). But before this can be done the operators have to be specified exactly.

For our discussion we always assume  $A = (a_{ij})_{i,j=1,\dots,n}$  to be symmetric and positive definite (s.p.d.) M-matrix. Accordingly we assume  $a_{ii} > 0$  and  $\sum_j a_{ij} \geq 0$  for all  $i = 1, \dots, n$ . On the interpolation operator  $I_H^h$  we assume full rank and the restriction operator is defined by

$$I_h^H := (I_H^h)^T \quad (98)$$

Consequently we define the coarse grid operator by

$$A^H := I_h^H A^h I_H^h \quad (99)$$

These type of setting is very convenient and is often referred to as Galerkin type. Because  $A_1$  is assumed to be s.p.d. all following  $A_q$  are s.p.d. (will be shown under sufficient preliminaries later) and consequently the coarse grid correction operators

$$T_q = I_q - I_{q+1}^q (A_{q+1}^{-1} I_{q+1}^q A_q) \quad (100)$$

become orthogonal projectors. The connection between the TG-operator of Point 5.1 and (100) is

$$M_q = S_q^{post} T_q S_q^{pre} \quad (101)$$

For the smoothing operator we would like to take the simplest one, namely the Gauß-Seidel or Jacobi relaxation. Of course, we could also take some block-smoothers, but this is not necessary (see [5]).

As we know, in a geometric multigrid there are various ways to choose the restriction and the coarse grid operator. In a pure AMG setting there hardly is another way to the choice above. Summarizing, we see that we just have to

define recursively the interpolation operator in an AMG algorithm.

Because the underlying problem is at most a discretization of an elliptic PDE of second order we would like to adopt what we mean by a grid. In geometrical MG the meaning is obvious; not so in AMG, where only the matrix and right hand side is available.

In the following we introduce ‘grids’ as a set of unknowns, which is assumed to be nested, i.e.

$$\omega_1 \supset \omega_2 \supset \dots \supset \omega_l \quad (102)$$

where  $\omega_i$  is the ‘grid’ on the  $i^{\text{th}}$  level. As we see in an AMG-like setting it is convenient to enumerate the levels in the opposite direction than in geometrical MG.

In terms of two consecutive levels  $h, H$  we assume that each coarse level variable  $u_k^H \in \omega_H$  is used to directly correct a uniquely defined fine grid variable  $u_{i(k)}^h \in \omega_h$ . Thus we can split up the set  $\omega_h$  in two disjoint subsets: the first contains the variables also presented in the coarser levels (C-variables or C-points) and the second one is just the complementary set (F-variables or F-points). Obviously we get

$$\sum_{j \in \omega_h} a_{ij}^h \underline{u}_j^h = \underline{f}_i^h \quad \forall i \in \omega_h \quad (103)$$

and after an appropriate renumbering of the C-points we get

$$\sum_{j \in C^h} a_{ij}^H \underline{u}_j^H = \underline{f}_i^H \quad \forall i \in C^h \quad (104)$$

By identifying each  $i \in \omega_h = C^h \cup F^h$  with a fictitious point in the plane we interpret the equations

$$A^h \underline{u}^h = \underline{f}^h \quad (105)$$

$$A^H \underline{u}^H = \underline{f}^H \quad (106)$$

as equations on the fictitious grid  $\omega_h$  and  $\omega_H$ , respectively, where (103) and (105) is nothing else then (97). This interpretation leads immediately to the following  $h$  to  $H$  coarsening process.

- define a splitting  $\omega_h = F^h \cup C^h$
- with  $\omega_H = C^h$  define interpolation weights  $\alpha_{ik}$  with the property

$$\alpha_{ik} = \delta_{ik} \quad \text{if } i \in C^h \quad (107)$$

and the interpolation can be represented by

$$(I_H^h e^H)_i = \sum_{k \in C^h} \alpha_{ik} e_k^H \quad \forall i \in \omega_h \quad (108)$$

- determine the restriction operator and the coarse grid operator by (98) and (99), respectively.

**Remark 5.9.**

- (i) For reference we call an interpolation of standard type, if it can be represented by (107) and (108).
- (ii) Obviously an interpolation of standard type has full rank.
- (iii) In practice the determination of a coarser grid and the computation of the interpolation weights are usually a very closely related process. Thus every time we talk about the construction of an interpolation we actually mean both.

**5.2.2 General convergence theory for the V-cycle**

Before we start to give some general results for the V-cycle, we give a definition for inner products and corresponding norms, which are very helpful in the following theory.

**Definition 5.10.** Let  $u, v \in \mathbb{R}^n$ ,  $A, D \in \mathbb{R}_n^n$  and  $D$  be a diagonal matrix with  $d_{ii} = a_{ii}$  for all  $i = 1, \dots, n$ . Then we define

$$\langle u, v \rangle_0 = \langle Du, v \rangle \tag{109}$$

$$\langle u, v \rangle_1 = \langle Au, v \rangle \tag{110}$$

$$\langle u, v \rangle_2 = \langle D^{-1}Au, Av \rangle \tag{111}$$

Here  $\langle \cdot, \cdot \rangle$  denotes the normal inner product in  $\mathbb{R}^n$  and  $\|\cdot\|_i$  the corresponding norm to  $\langle \cdot, \cdot \rangle_i$ .

First we would like to motivate the two major points in MG-method theory, i.e. the approximation assumption and the smoothing assumption.

**Theorem 5.11.** Let  $A$  be s.p.d. and assume the interpolation operator  $I_{q+1}^q$  to have full rank. The restriction is defined by (98). Furthermore suppose that for all  $e^q \in \mathbb{R}^q$

$$\|S_q e^q\|_1^2 \leq \|e^q\|_1^2 - \delta_1 \|T_q e^q\|_1^2 \tag{112}$$

$$\|S_q e^q\|_1^2 \leq \|e^q\|_1^2 - \delta_2 \|T_q S_q e^q\|_1^2 \tag{113}$$

hold with some  $\delta_1, \delta_2 > 0$  independent of  $e^q$  and  $q$ .

Then the following hold

- (i)  $\delta_1 < 1$  and provided that the coarsest grid is solved and that at least one smoothing step is performed after each coarse grid correction step, the V-cycle convergence factor is bounded from above by  $\sqrt{1 - \delta_1}$ .
- (ii) If at least one smoothing step is performed before each coarse grid correction step, then the V-cycle convergence factor is bounded from above by  $\sqrt{1/(1 + \delta_2)}$ .
- (iii) If at least one smoothing step is performed before and after each coarse grid correction step, then the V-cycle convergence factor is bounded from above by  $\sqrt{(1 - \delta_1)/(1 + \delta_2)}$ .

*Proof.* can be found in [28] □

**Theorem 5.12.** *Let  $\delta_1 = \alpha_1/\beta_1$  and  $\delta_2 = \alpha_2/\beta_2$ . Then*

$$\|S_q e^q\|_1^2 \leq \|e^q\|_1^2 - \alpha_1 \|e^q\|_2^2 \quad (114)$$

$$\|T_q e^q\|_1^2 \leq \beta_1 \|e^q\|_2^2 \quad (115)$$

*imply (112). Similarly*

$$\|S_q e^q\|_1^2 \leq \|e^q\|_1^2 - \alpha_2 \|S_q e^q\|_2^2 \quad (116)$$

$$\|T_q e^q\|_1^2 \leq \beta_2 \|e^q\|_2^2 \quad (117)$$

*imply (113).*

*Proof.* can be found in [28] □

**Remark 5.13.** *The two separated inequalities for  $S_q$  and  $T_q$  are called the smoothing assumption and the approximation assumption, respectively.*

How we can interpret these conditions: The first assumption in Theorem 5.11 means that the error components  $e^q$  that can not be effectively reduced by  $T_q$ , i.e.  $\|T_q e^q\|_1 \sim \|e^q\|_1$ , have to be effectively and uniformly reducible by  $S_q$ . Vice versa the error components  $e^q$  that can be effectively reduced by  $T_q$ , i.e.  $\|T_q e^q\|_1 \ll \|e^q\|_1$ , the smoothing operator  $S_q$  is allowed to be ineffective. Summarizing, these smooth errors, i.e.  $S_q e^q \sim e^q$ , have to be approximately in the range of the interpolation operator  $I_{q+1}^q$ . This will be the main objective to construct a good interpolation operator. What is left is a handy characterization of smooth errors.

### 5.2.3 The smoothing property

For positive definite matrices the Gauß-Seidel relaxation is known as effective smoother. Therefore, the next result will help us.

**Theorem 5.14.** *Let  $A \in \mathbb{R}_n^n$  as assumed. Let us define  $E := \sum_i |e_i|^2$  and  $R := \sum_i (|r_i|^2 / \sum_j |a_{ij}|^2)$ . Furthermore suppose that  $E^0, R^0$  and  $E^1, R^1$  be the values  $E, R$  before and after one Gauß-Seidel relaxation sweep, respectively, then*

$$E^1 \leq E^0 - \gamma_0 R^0 \quad (118)$$

$$E^1 \leq E^0 - \gamma_1 R^1 \quad (119)$$

where

$$\gamma_0 = (1 + \gamma_-)(1 + \gamma_+) - 1 \quad (120)$$

$$\gamma_1 = \gamma_- \gamma_+ - 1 \quad (121)$$

$$\gamma_- = \max_i \frac{\sum_{j < i} |a_{ij}|}{a_{ii}} \quad (122)$$

$$\gamma_+ = \max_i \frac{\sum_{j > i} |a_{ij}|}{a_{ii}} \quad (123)$$

*Proof.* can be found in [4, 5] □

The theorem shows, that fast convergence by a point relaxation scheme is always obtainable as long as  $R$  is comparable to  $E$ . The converse is also true; namely when  $R \ll E$ , then no point relaxation scheme can yield fast convergence. Summarizing we get

**Corollary 5.15.** *The convergence of a Gauß-Seidel relaxation scheme slow down, if and only if the normalized residuals  $R$  are small compared with the errors  $E$ .*

The last corollary gives rise to

**Definition 5.16.** *Let  $e \in \mathbb{R}^n$ . We call an error to be algebraically smooth if*

$$\|Se\|_1 \sim \|e\|_1 \quad (124)$$

**Corollary 5.17.** *In a straight forward manner we are able to present a practical characterization.*

- 1) *A necessary condition for an algebraically smooth error is  $\|e\|_2 \ll \|e\|_1$*
- 2) *If  $\|e\|_2 \ll \|e\|_1$  holds, then  $\|e\|_1 \ll \|e\|_0$*
- 3) *If  $\|e\|_2 \ll \|e\|_1$  holds, then  $\|e\|_2 \ll \|e\|_0$*

*Proof.* 1) By assumption the normalized residuals are small compared to the errors, thus the convergence of a Gauß-Seidel sweep slow down. Therefore we attain by definition a algebraically smooth error.

2) First we observe

$$\begin{aligned} \|e\|_1^2 &= \langle Ae, e \rangle \\ &= \langle D^{-1/2} Ae, D^{1/2} e \rangle \\ &\leq \sqrt{\langle D^{-1/2} Ae D^{-1/2} Ae, \rangle} \sqrt{\langle D^{1/2} e, D^{1/2} e \rangle} \\ &= \|e\|_2 \|e\|_0 \end{aligned} \quad (125)$$

Straight forward we get

$$\|e\|_1^2 \leq \|e\|_0 \|e\|_2 \ll \|e\|_0 \|e\|_1 \quad (126)$$

3) Follows immediately from 2). □

**Corollary 5.18.** *Let  $A \in \mathbb{R}_n^n$  be an M-matrix s.p.d., then we get the following characterization for an algebraically smooth error*

$$\sum_{i \neq j} \frac{|a_{ij}|}{a_{ii}} \frac{(e_i - e_j)^2}{e_i^2} \ll 1 \quad (127)$$

*Proof.* First we observe that

$$\langle Ae, e \rangle = 0.5 \sum_i \sum_j -a_{ij} (e_i - e_j)^2 + \sum_i (\sum_j a_{ij}) e_i^2 \quad (128)$$

From the condition for an algebraically smooth error  $\|e\|_1 \ll \|e\|_0$  we get

$$\sum_i \sum_j -a_{ij} (e_i - e_j)^2 + \sum_i (\sum_j a_{ij}) e_i^2 \ll \sum_i a_{ii} e_i^2 \quad (129)$$

Thus we can expect at least on the average

$$0.5 \sum_j -a_{ij} (e_i - e_j)^2 \ll a_{ii} e_i^2 \quad \forall i \quad (130)$$

and therefore

$$0.5 \sum_j \frac{-a_{ij} (e_i - e_j)^2}{a_{ii} e_i^2} \ll 1 \quad \forall i \quad (131)$$

□

The essence of the Corollary 5.18 is that a smooth error component varies generally slowly in the direction of strong coupling, i.e. from  $e_i$  to  $e_j$  if  $\frac{|a_{ij}|}{a_{ii}}$  is relatively large. This motivates to

**Definition 5.19.** *The neighborhood of a point  $i$  is defined by*

$$N(i) := \{j : a_{ij} \neq 0\} \quad (132)$$

*The distance between a fictitious point  $i$  and a set of fictitious points  $I$  is defined as*

$$d(i, I) := \frac{\sum_{j \in I} |a_{ij}|}{\max_{k \in I} |a_{ik}|} \quad (133)$$

*The set of strong coupled neighbors to  $i$  is defined as*

$$S^i := \{j \in N_i : d(i, \{j\}) \geq \alpha\} \quad (134)$$

*The set of points to which the point  $i$  has a strong coupling is defined as*

$$S^{i,T} := \{j \in N_i : i \in S^j\} \quad (135)$$

Concluding the subsection we are going to discuss a famous smoother, i.e. the GS smoother, and give sufficient conditions to fulfill the smoothing property.

**Corollary 5.20.** *Let  $A \in \mathbb{R}_n^n$ ,  $\gamma_-$  and  $\gamma_+$  as in (122) and (123).*

*Then the Gauß-Seidel relaxation satisfies (114) and (116) if*

$$\alpha_1 \leq \frac{1}{(1 + \gamma_-)(1 + \gamma_+)} \quad (136)$$

$$\alpha_2 \leq \frac{1}{\gamma_- \gamma_+} \quad (137)$$

*respectively.*

*Proof.* can be found in [28] □

**Remark 5.21.**

(i) For our class of matrices and if we apply a Gauß-Seidel forward/backward relaxation scheme we can calculate for  $\alpha_1$  and  $\alpha_2$  the values 0.25 and 1, respectively.

(ii) A similar result is also obtainable for the Jacobi smoother see [28].

**5.2.4 The approximation property**

**Two level convergence** Next we would like to discuss the approximation property and therefore we need an inequality of the form (approximation property)

$$\|T_h e^h\|_1^2 \leq \beta \|T_h e^h\|_2^2 \tag{138}$$

for all  $e^h \in \text{range}(T_h)$ . A more helpful characterization, which is sufficient for (138) provides the following theorem.

**Theorem 5.22.** Let  $A \in \mathbb{R}_n^n$  and let  $S_h$  satisfy the smoothing property. Suppose that the interpolation  $I_H^h$  has full rank and for every  $e^h \in \mathbb{R}^n$  the inequality

$$\min_{e^H} \|e^h - I_H^h e^H\|_0^2 \leq \beta \|e^h\|_1^2 \tag{139}$$

holds for some  $\beta > 0$  independent of  $e^h$ .

Then  $\beta \geq \alpha_1$  and the two level convergence factor satisfies

$$\|S_h T_h\|_1 \leq \sqrt{1 - \alpha_1 / \beta} \tag{140}$$

*Proof.* can be found in [28] □

We will give a characterization for an interpolation of standard form which relates the interpolations weights  $\alpha_{ij}$  with the matrix entries  $a_{ij}$ .

**Theorem 5.23.** Let  $A \in \mathbb{R}_n^n$ . Furthermore suppose that for any subset  $C^h$  of  $C$ -points the interpolation  $I_H^h$  is of the standard form with  $\alpha_{ik} \geq 0$  and  $s_i \leq 1$ , where  $s_i := \sum_{k \in C} \alpha_{ik}$ .

Then the property (139) is fulfilled if the two inequalities hold with  $\beta > 0$  independent of  $e = e^h \in \mathbb{R}^n$ .

$$\sum_{i \in F} \sum_{k \in C} a_{ii} \alpha_{ik} (e_i - e_k)^2 \leq \frac{\beta}{2} \sum_{i,j} (-a_{ij}) (e_i - e_j)^2 \tag{141}$$

$$\sum_{i \in F} a_{ii} (1 - s_i) e_i^2 \leq \beta \sum_i \left( \sum_j a_{ij} \right) e_i^2 \tag{142}$$

*Proof.* can be found in [28] □

Concluding we have shown the convergence for a two level method; unfortunately we can not give a proof for a multilevel method, but we will give some heuristics to obtain a good multilevel cycle.



**Remarks on the multilevel method** To set up a multilevel cycle we need the operator on the coarse levels to have the same properties as on the finest grid. Therefore we give a Theorem which ensures us this condition for M-matrices.

**Theorem 5.24.** *Let  $A^h \in \mathbb{R}_n^n$  be a symmetric, weakly diagonally dominant M-matrix and let the interpolation weights satisfy*

$$0 \leq a_{ii}\alpha_{ik} \leq \beta|a_{ik}| \quad \forall i \in F, k \in C^i \quad (143)$$

$$0 \leq a_{ii}(1 - s_i) \leq \beta \sum_j a_{ij} \quad \forall i \in F, k \in C^i \quad (144)$$

with some  $\beta \leq 2$ .

Then  $A^H$  is also a symmetric, weakly diagonally dominant M-matrix.

*Proof.* can be found in [28] □

As mentioned above, we can not overcome the q-dependency of the multilevel cycle. First of all we could overcome it by ‘better’ cycles but this is not satisfactory.

The problem is that we can not interpolate the smooth error well enough. To get a q-independent V-cycle, we have to improve the approximation property. As we know, a stronger condition for the approximation property is

$$\|T_h e^h\|_1^2 \leq \beta \|e^h\|_2^2 \quad (145)$$

Consequently we have to connect this condition with Theorem 5.22 to get a two level convergence rate first.

**Lemma 5.25.** *If the inequality (145) holds, then we also have*

$$\min_{e^H} \|e^h - I_H^h e^H\|_0^2 \leq \beta^2 \|e^h\|_2^2 \quad (146)$$

*Proof.* can be found in [28] □

Obviously this condition is much stronger for smooth errors, because  $\|e^h\|_2 \ll \|e^h\|_1$ . Accordingly the new condition (145) increases the order of interpolation.

**Remark 5.26.**

- (i) *Theoretically we could try to achieve an interpolation formula which satisfies the new condition rather than the old one. But unfortunately it is hardly possible just by using only algebraic information.*
- (ii) *In practice it turns out that we do not have to fulfill (145) exactly, we just need certain objectives (see Point 5.2.6).*

### 5.2.5 Interpolation formulas

As we saw in point 5.2.3 we can characterize a smooth error  $e$  by

$$\|e\|_2 \ll \|e\|_1 \quad (147)$$

or in an explicit way

$$\sum_i \frac{r_i^2}{a_{ii}} \ll \sum_i r_i e_i \quad \forall i \in \omega_h \quad (148)$$

Therefore we can expect (at least on the average)

$$|r_i| \ll a_{ii}|e_i| \quad \forall i \in \omega_h \quad (149)$$

On the other hand we know

$$r_i = a_{ii}e_i + \sum_{j \in N_i} a_{ij}e_j \quad \forall i \in \omega_h \quad (150)$$

As mentioned earlier, we know that any error that is slow to convergence should be well presented by the range of interpolation. Hence if we combine (149) and (150) we get immediately for smooth errors

$$a_{ii}e_i + \sum_{j \in N_i} a_{ij}e_j \sim 0 \quad \forall i \in \omega_h \quad (151)$$

For reasons of simplicity let us consider an underlying grid  $\Omega_h$  and a matrix arising from an FE-discretization. For practical reasons we interpret (152) as

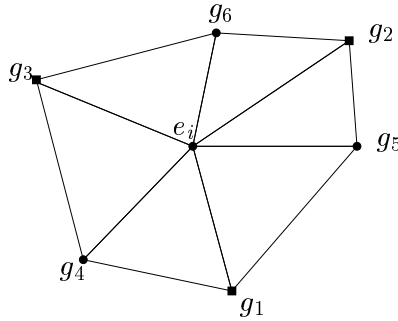


Figure 10: A part of a grid  $\omega_h$

$$a_{ii}e_i + \sum_{j \in N_i} a_{ij}e_j = 0 \quad \forall i \in \omega_h \quad (152)$$

If we look at Figure 10 and at equation (152) we see this is nothing but the harmonic extension with the polygonal-line of the  $N_i$  as boundary. Further on we know the values in the points  $C^i = N_i \cap S^i$  and it is clear that  $g_j$  for all  $j \in F^i$  are not determined yet. These points can be calculated in two ways. First we determine them by simple linear interpolation, i.e. we interpolate a

point  $j \in F^i$  with the neighboring coarse grid points. With a straightforward calculation we get the following interpolation weights

$$\tilde{\alpha}_{ij} = \begin{cases} 1 & i = j \in C \\ -\frac{a_{ij} + \tilde{c}_{ij}}{a_{ii}} & i \in F, j \in C^i \\ 0 & \text{else} \end{cases} \quad (153)$$

where

$$\tilde{c}_{ij} = \sum_{k \in F^i} \frac{a_{ik} a_{kj}}{\sum_{l \in C^i} a_{kl}} \quad (154)$$

The other way to determine the interpolation weights, which behaves better in numerical situations, is to determine the fine grid points, i.e.  $j \in F^i$ , in a more sophisticated way without losing the sparsity of the matrix  $A_h$  to  $A_H$ . In this way we get

$$\alpha_{ij} = \begin{cases} 1 & i = j \in C \\ -\frac{a_{ij} + c_{ij}}{a_{ii} + c_{ii}} & i \in F, j \in C^i \\ 0 & \text{else} \end{cases} \quad (155)$$

where

$$c_{ij} = \sum_{k \in F^i} \frac{a_{ik} a_{kj}}{\sum_{l \in C^i} a_{kl} + a_{ki}} \quad (156)$$

Heuristically we could state that the last interpolation is better. This is crucial because the better the interpolation, the better the MG-rate. Both derivations of the interpolation weights  $\tilde{\alpha}_{ij}$  and  $\alpha_{ij}$  can be found in [17].

### 5.2.6 The coarsening algorithm and remarks for implementation

**General remarks and notations** For the following algorithms we denote the system matrix, the right hand side and the solution by  $A_h$ ,  $\underline{f}_h$  and  $\underline{u}_h$ , respectively. The abbreviations  $KA$  and  $WA$  stand for an array of the system and the interpolation matrices, respectively. We always identify the system matrix on the finest level with the number one and the coarser levels with  $q$ . In the Algorithm 2 a possible structure for an interface is given. The procedure **Solver** can be a normal RSAMG cycle or a preconditioned CG.

**The coarsening algorithm** As mentioned, for RSAMG we have to construct the coarse grid matrices purely algebraically or at least with less information about the underlying grid. This is done by constructing an interpolation operator which depends on the matrix entries and then using Galerkins method. An algorithm which supplies the coarse grid matrices and the interpolation operator see Algorithm 3.

The efficiency of the coarsening process should be the basis of an efficient solution process of a multigrid cycling type. On the one hand we can say

```

PROCEDURE Main()
BEGIN

    get the system matrix  $A_h$ 
    get the right hand side  $\underline{f}_h$ 
     $KA_1 = A_h$ 
    Setup( $KA, WA, 1$ )
    Solver( $KA, WA, \underline{f}_h, \underline{u}_h$ )
END

```

Algorithm 2: Algorithm for positive definite matrices

that the h-level convergence factor is better the more coarse grid points we have on the coarse level. On the other hand the amount of work necessary for relaxation on the H-level increases with the size of the coarse grid matrix. Hence it is advantageous to limit the number of C-points. In order to obtain a good interpolation formula and to keep the numerical work small we suggest two criteria.

- (C1) For each  $i \in F$ , each point  $j \in S^i$  should either be in C, or should be strongly connected to at least one point in  $C^i$ .
- (C2) C should be a maximal subset of all points with the property that not two C-points are strongly connected to each other.

However, there are examples where it is impossible to fulfill both criteria, but in practice they perform well.

In order to keep the work involved in the choice of the coarse grid small, a two part process is presented. In the first part of the process we select coarse grid points in such a manner, that as many strong couplings as possible are taken into account. The second part of the process has to make sure, that each point in F is strongly coupled with points in C.

### 5.2.7 The algorithm for general s.p.d. matrices

Constructive to the described algorithm we make a generalization to all matrices of the s.p.d. class. This is of course very important because the M-matrix property can be lost if

- we use quadratic ansatz-functions instead of linear ones
- the triangulation is not regular
- we use long thin rectangles

In any of these cases we will not lose our symmetry and positive definiteness. But if the non-negative values in the matrix get too big, the described algorithm will break down. Especially we get negative interpolation weights. How can we remove this leak?

```

PROCEDURE Setup( $KA, WA, level$ )
BEGIN

  IF coarse grid is small enough THEN
    BEGIN
      RETURN
    END
  ELSE
    BEGIN
      (* choose coarse grid *)
       $K = KA_{level}$ 
       $W = WA_{level}$ 
      (* prepare the strong couplings *)
      Setup0( $K, S^i, S^{i,T}$ )
      (* compute the coarse and fine grid points *)
      Setup1( $K, C, F$ )
      Setup2( $C, F$ )
      (* define the interpolation and restriction operator *)
      InterpolationWeights( $K, W$ )
      (* compute the coarse grid matrix *)
      CoarseMatrix( $K, W, KA_{level+1}$ )
      Setup( $KA, WA, level + 1$ )

    END

  END

```

Algorithm 3: Principle setup phase

Notice that AMG could be interpreted as a preconditioner (see point 5.1) and therefore we just need a spectral equivalent matrix on the finest grid. Hence we are looking for a matrix  $B_l$  with

$$B_l \sim K_l \tag{157}$$

and  $B_l$  is M-matrix. A very sophisticated method to get a spectrally equivalent M-matrix is to extract more information out of the element-matrices.

**Theorem 5.27.** *Let  $A \in \mathbb{R}_n^n$  be a s.p.d. matrix arising from an FE discretization. Moreover let  $A^{(r)}$  be the corresponding element matrices and assume  $\tilde{A}^{(r)} \sim A^{(r)}$  for all  $r \in \mathbb{R}_h$ . (Notice:  $\tilde{A}^{(r)}$  and  $A^{(r)}$  are positiv semidefinite and thus  $N(\tilde{A}^{(r)}) = N(A^{(r)})$ .) Then we get*

$$A \sim \tilde{A} \tag{158}$$

where  $\tilde{A}$  is the assembled matrix of the  $\tilde{A}^{(r)}$ .

```

PROCEDURE Setup1( $A, C, F$ )
BEGIN

   $C = \emptyset, F = \emptyset$ 
  WHILE  $C \cup F \neq |A|$  DO
  BEGIN
    Pick( $i \in N \setminus (C \cup F)$ )
    IF  $|S^{i,T}| + |S^{i,T} \cap F|$  THEN
    BEGIN
       $F = N \setminus C$ 
    END
    ELSE
    BEGIN
       $C = C \cup i$ 
       $F = F \cup (S^{i,T} \setminus C)$ 
    END
  END

END

```

Algorithm 4: Setup phase I

*Proof.* Let us suppose  $u \in \mathbb{R}^n$ , then we attain for the ‘energy’

$$\begin{aligned}
u^T A u &= u^T \left( \sum_{r \in \mathbb{R}_h} C_r A^{(r)} C_r^T \right) u \\
&= \sum_{r \in \mathbb{R}_h} u^T C_r A^{(r)} (u^T C_r)^T
\end{aligned} \tag{159}$$

where  $C_r$  are the connectivity matrices. Consequently we get

$$\begin{aligned}
u^T A u &\leq \sum_{r \in \mathbb{R}_h} \bar{\gamma}^{(r)} u^T C_r \tilde{A}^{(r)} (u^T C_r)^T \\
&\leq \max_{r \in \mathbb{R}_h} \left\{ \bar{\gamma}^{(r)} \right\} u^T \tilde{A} u
\end{aligned} \tag{160}$$

$$\begin{aligned}
u^T A u &\geq \sum_{r \in \mathbb{R}_h} \underline{\gamma}^{(r)} u^T C_r \tilde{A}^{(r)} (u^T C_r)^T \\
&\geq \min_{r \in \mathbb{R}_h} \left\{ \underline{\gamma}^{(r)} \right\} u^T \tilde{A} u
\end{aligned} \tag{161}$$

Summarized we have shown

$$A^{(r)} \sim \tilde{A}^{(r)} \quad \forall r \in \mathbb{R}_h \Rightarrow A \sim \tilde{A} \tag{162}$$

□

To use the above theorem we have to determine for each element  $\delta_r$ ,  $r \in \mathbb{R}_h$ , a spectrally equivalent element matrix; but how could this be done?

From a practical point of view we have to check the element-matrix before we assemble it and if in the element-matrix  $A^{(r)}$  is no M-matrix, we simply replace  $A^{(r)}$  by  $\tilde{A}^{(r)}$ .

In the following discussion we restrict ourselves to  $\Omega \subset \mathbb{R}^2$  and the potential equation, i.e.

$$-\operatorname{div}(\mu(x) \operatorname{grad} u(x)) = f(x) \quad x \in \Omega \quad (163)$$

Next we consider long thin rectangles as in Figure 11 with bilinear ansatz-functions. By calculating the element stiffness matrix we attain

$$A^{(r)} = \frac{1}{6pq} \begin{pmatrix} 2(p^2 + q^2) & p^2 - 2q^2 & -2p^2 + q^2 & -(p^2 + q^2) \\ p^2 - 2q^2 & 2(p^2 + q^2) & -(p^2 + q^2) & -2p^2 + q^2 \\ -2p^2 + q^2 & -(p^2 + q^2) & 2(p^2 + q^2) & p^2 - 2q^2 \\ -(p^2 + q^2) & -2p^2 + q^2 & p^2 - 2q^2 & 2(p^2 + q^2) \end{pmatrix} \quad (164)$$

Notice, that for  $q < \frac{p}{\sqrt{2}}$  positiv off-diagonal elements will occur and thus the M-matrix property is lost. Immediately we get three objectives for  $\tilde{A}^{(r)}$ :

- $\tilde{A}^{(r)}$  is an M-matrix
- $\tilde{A}^{(r)}$  is algebraically computeable from  $A^{(r)}$
- $\tilde{A}^{(r)\dagger} A^{(r)}$  has a constant spectrum, independent from ‘bad’ parameters

where  $\tilde{A}^{(r)\dagger}$  is the pseudo-inverse of  $\tilde{A}^{(r)}$ . One method to get a spectrally equivalent M-matrix  $\tilde{A}^{(r)}$  to  $A^{(r)}$  is the following:

- cut the rectangle into two triangles (see Figure 11)
- assume on the triangles linear ansatz-functions and compute the element matrices on the triangles
- assemble the two element matrices

In this way we obtain the matrix

$$\tilde{A}^r = \frac{1}{2pq} \begin{pmatrix} p^2 + q^2 & -q^2 & -p^2 & 0 \\ -q^2 & p^2 + q^2 & 0 & -p^2 \\ -p^2 & 0 & p^2 + q^2 & -q^2 \\ 0 & -p^2 & -q^2 & p^2 + q^2 \end{pmatrix} \quad (165)$$

Now it is easy verified that (165) is symmetric, positiv semidefinite and is an M-matrix. Another method to get a spectrally equivalent M-matrix is the Algorithm 7. Therefore let us assume the numbering of the unknowns as in Figure 11 and the original element stiffness matrix (164). Consequently we get for our element matrix  $A^{(r)}$  the spectrally equivalent matrix with Algorithm 7.

$$\tilde{A}^r = \frac{1}{6pq} \begin{pmatrix} 2(p^2 + q^2) & -3q^2 & -2p^2 + q^2 & 0 \\ -3q^2 & 2(p^2 + q^2) & 0 & -2p^2 + q^2 \\ -2p^2 + q^2 & 0 & 2(p^2 + q^2) & -3q^2 \\ 0 & -2p^2 + q^2 & -3q^2 & 2(p^2 + q^2) \end{pmatrix} \quad (166)$$

Again this matrix is symmetric, positiv semidefinite and an M-matrix. Finally it is shown that (165) and (166) are spectrally equivalent to (164).

**Theorem 5.28.** *Let (164) be an element matrix. Then (165) and (166) are spectrally equivalent matrices to (164). The eigenvalues of the generalized eigenvalue problem are  $\{0, 1, 1, \frac{1}{3}\}$  and  $\{0, \frac{1}{2}, 1, \frac{3p^2}{2p^2-q^2}\}$  for (165) and (166), respectively.*

*Proof.* Let us consider the generalized eigenvalue problem

$$A^{(r)}u = \lambda \tilde{A}^{(r)}u \quad (167)$$

Then by computing the eigenvalues for  $\tilde{A}^{(r)\dagger}A^{(r)}$  the desired result follows.  $\square$

**Remark 5.29.**

- (i) *Obviously the spectral constants for (166) depends on  $p$  and  $q$ . But as assumed  $p > q$  and therefore  $\frac{3p^2}{2p^2-q^2}$  is bounded.*
- (ii) *The above technique can also be applied to general quadrilaterals if the inner angels are bounded.*
- (iii) *For the first method  $p$  and  $q$  or at least the ratio  $\frac{p}{q}$  has to be known.*

For a triangle like in Figure 12 which is in general no M-matrix for  $0 < q < 1$  it is the best to cut it into 2 triangles like in Figure 12.

A technique to find the best possible M-matrix for a general element matrix  $A^{(r)}$  can be found in Algorithm 8. This procedure is very general and useful for all arising element matrices, which are positiv semidefinite. Notice, that this technique can be also applied in 3D.



**PROCEDURE Setup2**( $C, F$ )  
**BEGIN**

$T = \emptyset$   
    **WHILE**  $T \subset F$  **DO**  
    **BEGIN**  
        *pick*  $i \in F \setminus T$   
         $T = T \cup \{i\}$   
         $\tilde{C} = \emptyset$   
         $C^i = S^i \cap C$   
         $F^i = S^i \cap F$   
        **WHILE**  $F^i \neq \emptyset$  **DO**  
        **BEGIN**  
            *pick*  $j \in F^i$   
             $F^i = F^i \setminus \{j\}$   
            **IF**  $\frac{d(j, C^i)}{d(i, \{j\})} \leq \beta$  **THEN**  
            **BEGIN**  
                **IF**  $|C^i| = 0$  **THEN**  
                **BEGIN**  
                     $\tilde{C} = \{j\}$   
                     $C^i = C^i \cup \{j\}$   
                **END**  
                **ELSE**  
                **BEGIN**  
                     $C = C \cup \{i\}$   
                     $F = F \setminus \{i\}$   
                **END**  
            **END**  
        **END**  
         $C = C \cup \tilde{C}$   
         $F = F \setminus \tilde{C}$   
    **END**

**END**

Algorithm 5: Setup phase II

```

PROCEDURE ElementMatrixCheck( $A^{(r)}$ )
BEGIN

  (* get the element matrix for element r, i.e.  $A^{(r)}$  *)
  IF element matrix  $A^{(r)}$  is no M-matrix THEN
    BEGIN
      (* produce a spectral equivalent M-matrix *)
      SpectralMatrix $A^{(r)}$ 
    END
  Assemble( $A^{(r)}$ )
END

```

Algorithm 6: Algorithm for positive definite matrices

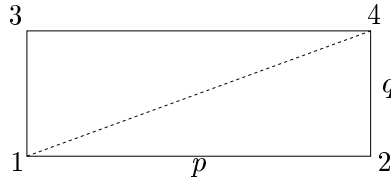


Figure 11: A long thin rectangle splitted into two triangles

```

PROCEDURE SpectralMatrix( $A^{(r)}$ )
BEGIN

  IF  $a_{12} > 0$  or  $a_{34} > 0$  THEN
    BEGIN
       $s_x = 2(a_{12} + a_{14} + a_{23} + a_{34})$ 
       $s_y = 2(a_{13} + a_{24})$ 
    END
  IF  $a_{13} > 0$  or  $a_{24} > 0$  THEN
    BEGIN
       $s_x = 2(a_{12} + a_{34})$ 
       $s_y = 2(a_{13} + a_{24} + a_{14} + a_{23})$ 
    END
  calculate the elements of  $\tilde{A}^{(r)}$ 
   $\tilde{a}_{11} = \tilde{a}_{22} = \tilde{a}_{33} = \tilde{a}_{44} = -\frac{s_x + s_y}{2}$ 
   $\tilde{a}_{12} = \tilde{a}_{34} = \frac{s_x}{2}$ 
   $\tilde{a}_{13} = \tilde{a}_{24} = \frac{s_y}{2}$ 
   $\tilde{a}_{14} = \tilde{a}_{23} = 0$ 
END

```

Algorithm 7: Algebraic technique to attain an M-matrix

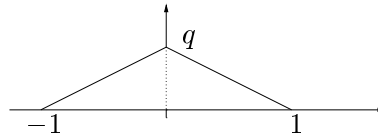


Figure 12: A ‘thin’ triangle splitted into two triangles

```

PROCEDURE GeneralSpectralMatrix( $A^{(r)}$ )
BEGIN

  get the element matrix for element r, i.e.  $A^{(r)}$ 
  IF element matrix  $A^{(r)}$  is no M-matrix THEN
    BEGIN
      calculate  $\tilde{A}^{(r)}$  with
       $\frac{\lambda_{max}}{\lambda_{min}} \rightarrow \min$ 
      such that  $\tilde{A}^{(r)}$  is an M-matrix, symmetric, positiv semidefinite
    END
  END

```

Algorithm 8: General technique to get an M-matrix

## 6 Numerical results

Up till now we have developed the theory to get an optimal solver. In that section we are going to show the effectivity and robustness of RSAMG. Therefore the code was implemented in the program package FEPP, see [31] and RSAMG is used as a preconditioner for CG.

The initial discretization was performed by hand and the refinement of the mesh was done uniformly by the program package FEPP. As examples we treat the potential equation with Dirichlet boundary conditions, which typically arises in magnetic field problems. Before, some CPU-time measurements are presented for the unit square.

All calculations were done on a SUN Ultra 1 workstation.

### 6.1 First example

Let us first consider the homogeneous Dirichlet boundary value problem for the Poisson equation, i.e.

$$-\Delta u = 1 \quad \forall x \in \Omega \quad (168)$$

$$u = 0 \quad \forall x \in \partial\Omega \quad (169)$$

where  $\Omega = (0, 1)^2$ . Further let  $\Omega$  be provided with a uniform triangulation for the FE discretization with linear ansatz functions. In Table 1 the CPU-time for the setup and the solution phase are shown.

unknowns	level	setup	solution	iterations
1089	6	1.15	0.12	5
4225	7	5.10	0.70	6
16641	8	21.08	3.02	6

Table 1: Demanded CPU-time in seconds

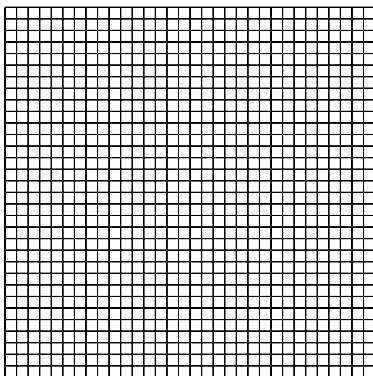


Figure 13: Initial matrix graph

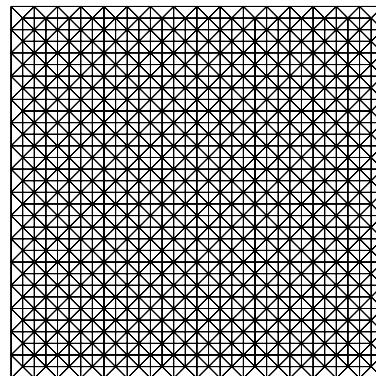


Figure 14: one coarsening step

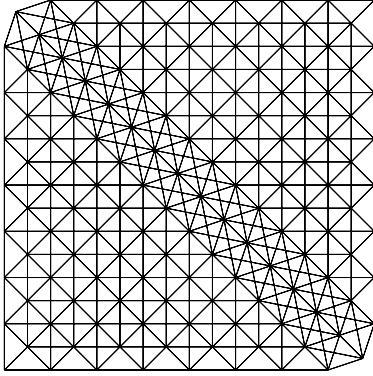


Figure 15: two coarsening steps

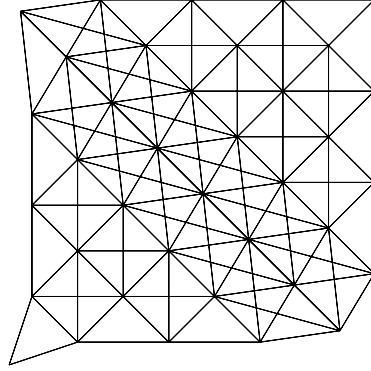


Figure 16: three coarsening steps

## 6.2 Shielding problem with one thin material

Now a classical shielding problem is considered. Therefore we assume a geometry as in Figure 17 and the equations

$$-\operatorname{div}(\mu(x) \operatorname{grad} u(x)) = 0 \quad \forall x \in \Omega_i \quad (170)$$

$$u(x) = -\mu_1(x) \cdot (B_x x + B_y y) \quad \forall x \in \partial\Omega \quad (171)$$

$$-\mu_1 \frac{\partial u_1(x)}{\partial n} = -\mu_2 \frac{\partial u_2(x)}{\partial n} \quad \forall x \in \partial\Omega_b \quad (172)$$

I.e. it is the system arising from the scalar potential in 2D (see Section 2). For the constants see Table 2.

example number	$\mu_1$	$\mu_2$	$B_x$	$B_y$
1	1	10e-6	21	45
2	1	10e-3	21	45
3	1	10e+3	21	45
4	1	10e+6	21	45

Table 2: Assumed constants for the shielding problem with one material

In the following examples we will use the abbreviations  $\kappa$ ,  $\rho$  and  $t$  for the condition number, the convergence factor and the thickness of the object, respectively.

First we have to emphasize that the element preconditioning technique is necessary to attain a robust solver. This can be seen in Table 3 where a result without element preconditioning is shown. In Table 4, 5, 6, 7 numerical examples are presented for different thickness of the body. For Examples 1 to 4 we can observe that the condition number  $\kappa = \kappa(C^{-1}A)$  remains constant. Unfortunately the number of iterations is slightly growing in the Examples 3 and 4.

unknowns	level	$\kappa$	$\rho$	t	iterations
2145	7	63	0.78	1e-2	19
		2563	0.96	1e-3	27
		218807	0.99	1e-4	36

Table 3: Example 4 without element preconditioning

unknowns	level	$\kappa$	$\rho$	t	iterations
153	5	1.85	0.15	1e-1	10
		1.79	0.14	1e-2	10
		1.79	0.14	1e-3	10
		1.79	0.14	1e-4	10
561	6	2.18	0.19	1e-1	11
		1.89	0.16	1e-2	11
		1.89	0.16	1e-3	11
		1.88	0.16	1e-4	11
2145	7	4.21	0.34	1e-1	13
		1.99	0.17	1e-2	12
		1.98	0.17	1e-3	12
		1.98	0.17	1e-4	12

Table 4: Example 1;  $\mu_1 = 1$ ,  $\mu_2 = 10^{-6}$

unknowns	level	$\kappa$	$\rho$	t	iterations
153	5	1.85	0.15	1e-1	10
		1.75	0.14	1e-2	10
		1.70	0.13	1e-3	11
		2.46	0.22	1e-4	12
561	6	2.17	0.19	1e-1	11
		1.82	0.15	1e-2	11
		1.82	0.15	1e-3	13
		2.13	0.19	1e-4	14
2145	7	4.21	0.34	1e-1	13
		1.97	0.18	1e-2	12
		1.89	0.16	1e-3	13
		1.88	0.16	1e-4	15

Table 5: Example 2;  $\mu_1 = 1$ ,  $\mu_2 = 10^{-3}$

unknowns	level	$\kappa$	$\rho$	t	iterations
153	5	4.71	0.37	1e-1	17
		4.45	0.37	1e-2	18
		4.02	0.33	1e-3	17
		3.28	0.29	1e-4	15
561	6	4.96	0.38	1e-1	21
		4.44	0.37	1e-2	20
		4.04	0.34	1e-3	20
		3.34	0.29	1e-4	18
2145	7	8.32	0.49	1e-1	27
		4.95	0.38	1e-2	24
		4.60	0.36	1e-3	23
		3.52	0.31	1e-4	20

Table 6: Example 3;  $\mu_1 = 1$ ,  $\mu_2 = 10^{+3}$

unknowns	level	$\kappa$	$\rho$	t	iterations
153	5	4.72	0.37	1e-1	16
		4.53	0.36	1e-2	16
		4.53	0.36	1e-3	16
		4.52	0.36	1e-4	16
561	6	4.97	0.38	1e-1	20
		4.50	0.36	1e-2	18
		4.50	0.36	1e-3	18
		4.49	0.36	1e-4	19
2145	7	8.28	0.48	1e-1	24
		4.99	0.38	1e-2	21
		4.96	0.38	1e-3	22
		4.96	0.38	1e-4	23

Table 7: Example 4;  $\mu_1 = 1$ ,  $\mu_2 = 10^{+6}$

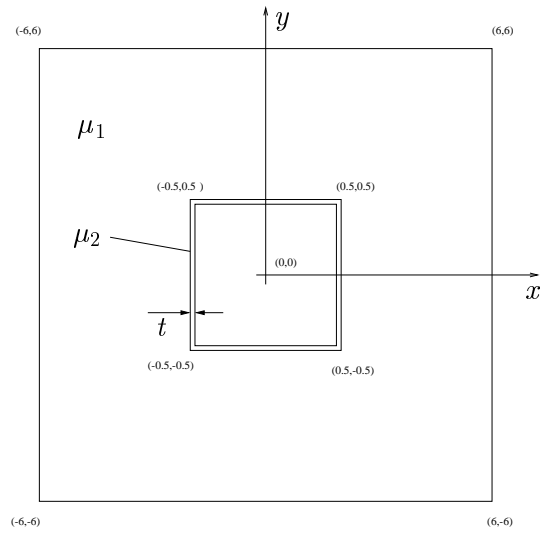


Figure 17: Example with one material

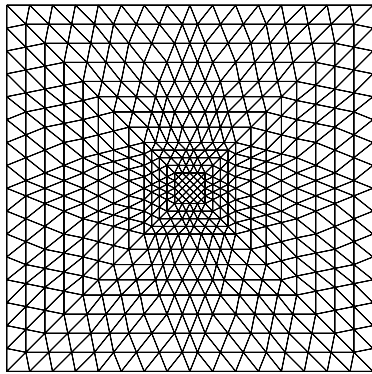


Figure 18: Initial matrix graph

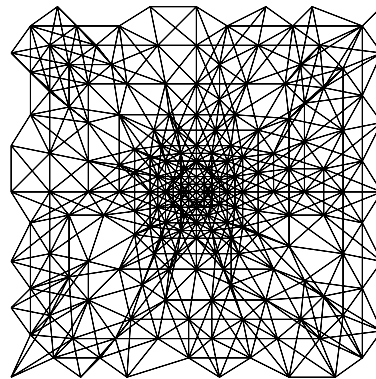


Figure 19: one coarsening step

### 6.3 Shielding problem with different material layers

Now we will compose a kind of sandwich material, see Figure 23. Again we solve the problem (170), (171), (172). For the parameters see Table 8.

Like in Example 3 and 4 we observe a slightly growing in the number of iterations; the condition number  $\kappa$  remains constant again, and therefore also the convergence factor  $\rho$ .



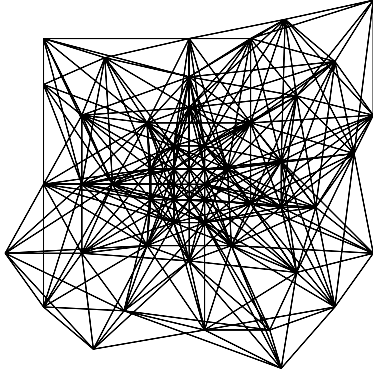


Figure 20: two coarsening steps

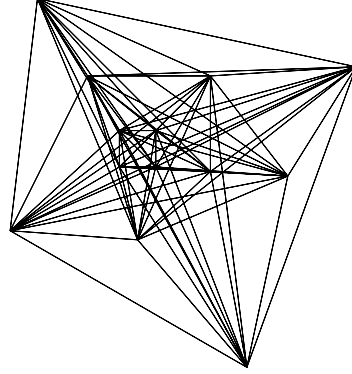


Figure 21: three coarsening steps

example number	$\mu_1$	$\mu_2$	$\mu_3$	$B_x$	$B_y$	$t_1$	$t_2$
5	1	10e+3	10e-3	21	45	10e-3	10e-2
6	1	10e+6	10e-6	21	45	10e-3	10e-2

Table 8: Assumed constants for the shielding problem with different materials

unknowns	level	$\kappa$	$\rho$	iterations
169	6	4.53	0.35	16
593	7	4.50	0.35	19
2209	8	4.96	0.37	21

Table 9: Example 5;  $\mu_1 = 1$ ,  $\mu_2 = 10^{+3}$ ,  $\mu_3 = 10^{-3}$

unknowns	level	$\kappa$	$\rho$	iterations
169	6	4.25	0.36	19
593	7	4.33	0.36	20
2209	8	4.80	0.38	24

Table 10: Example 6;  $\mu_1 = 1$ ,  $\mu_2 = 10^{+6}$ ,  $\mu_3 = 10^{-6}$

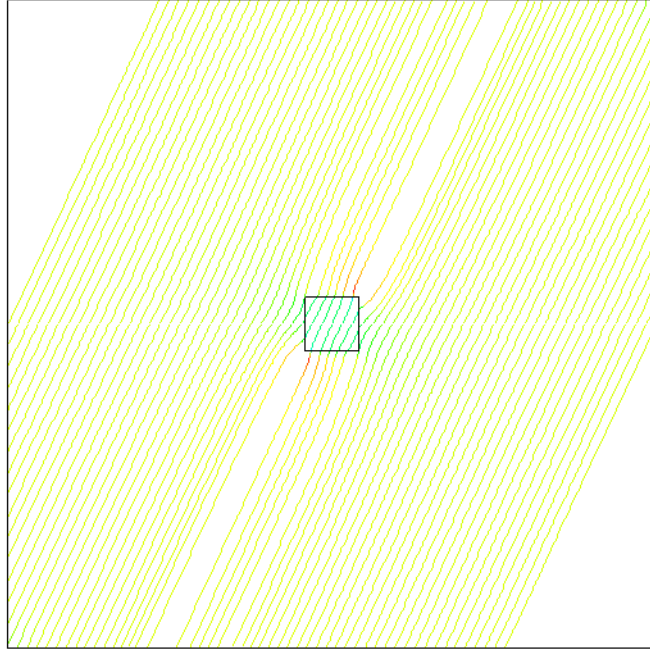


Figure 22: Streamlines of Example 1

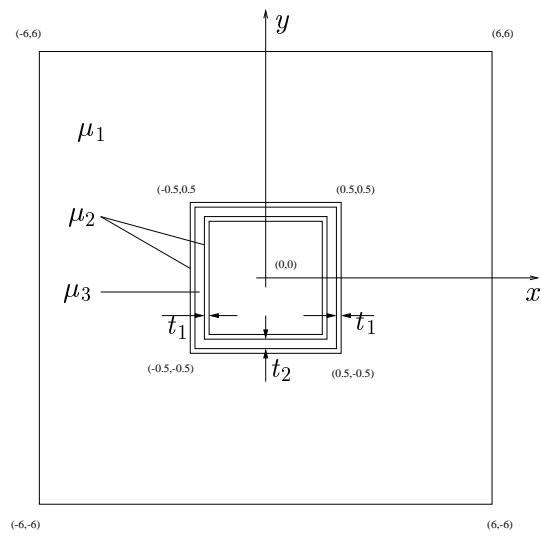


Figure 23: Example with two materials

## 7 Conclusions and further remarks

In the last sections a lot of results were presented to build a good model and solution strategy for a magnetic shielding problem. First of all we started with a short introduction, where the problem is heuristically mentioned and some literature on the arising topics is provided.

Further a practical problem and its physical background is posed. There the Maxwell equations for stationary objects are reduced to a very simple vector valued equation. Two mathematical methods were given to reduce at least in 2D the vector valued equation to a scalar PDE. It was shown that for the scalar potential the reduction to a scalar PDE works also well in the 3D case. Constructively a weak formulation was submitted for the scalar potential and in this context an existence and uniqueness theory.

In order to build up an FE-space we were employed to find a discretization of micro structures in 2D. Two possibilities were presented; first a method with long thin triangles and secondly a technique with long thin quadrilateral. Accordingly also refinement strategies were suggested.

Constructively an FE-space was build with underlying linear and bilinear form-functions. It was shown that for long thin triangles and rectangles the condition-number behaves like  $O(h_{min}^{-2})$ , where  $h_{min}$  denotes the global minimal element-edge. Further a convergence result was shown in  $H^1$ .

Consequently a numerical solution strategy of the arising finite, linear, sparse and bad conditioned system was discussed. Three different areas were discussed where RSAMG behaves favorable. After that a motivation and a construction of the RSAMG-method was given. In this context a two level convergence rate estimate was proved and some remarks on the multilevel case were provided. In order to apply the RSAMG-method to problems with with general symmetric positive definite matrices we suggested two techniques.

Finally we presented numerical examples where RSAMG had proven to be a robust and efficient black-box solver and preconditioner for CG at least in a numerical way. Unfortunately the setup phase is expensive, which is the main disadvantage of RSAMG against other solvers. Although the presented numerical examples were very academical and therefore ‘nice’, normal geometrical MG-methods would break down.

One thing has to be emphasized: RSAMG is a robust solver in a wide range of applications and especially well suited for complex domains. Therefore it makes it very attractive to structural optimization or even in fluid dynamics. Consequently, a generalization of RSAMG would be due to systems of PDEs. Other areas of application for RSAMG could be nonlinear problems as solver in a Newton-method, in time-dependent problems for the elliptic part or in purely discrete problems.

Actually the last comments were only in the sense of generalizing the algorithm. A further development of the model could also be faced.

## References

- [1] R.A.Adams  
*Sobolev Spaces*, Academic Press, New York, 1975
- [2] T.Apel, M.Dobrowolski  
*Anisotropic Interpolation with Applications to the Finite Element Method*,  
Computing 47, 1992
- [3] J.H.Bramble, J.E.Pasciak, J.Xu  
*Parallel multilevel preconditioners*, Math. Comput. 55, 1990
- [4] A.Brandt  
*Algebraic Multigrid Theory: The Symmetric Case*, in Preliminary Proceedings of the International Multigrid Conference, Copper Mountain, Colorado, (S.McCormick and U.Trottenberg ed),1983
- [5] A.Brandt, S.McCormick, J.W.Ruge  
*Algebraic Multigrid (AMG) for sparse matrix equations*, in Sparsity and it's Application, (D.J.Evans ed), 1984
- [6] P.G. Ciarlet  
*The Finite Element Method for Elliptic Problems*, North-Holland Publishing Company, Amsterdam-New York-Oxford, 1987
- [7] J.Dendy  
*Blackbox multigrid for nonsymmetric problems*, Appl.Math.Comput.,1983
- [8] H.W. Engl  
*Partielle Differentialgleichungen I*, Skriptum zur gleichnamigen Vorlesung, Johannes Kepler Universität, Ordinariat für Industriemathematik, Linz, 1995
- [9] J.Fuhrmann  
*A Modular Algebraic Multilevel Method*, Weierstraß Institute for Applied Analysis and Stochastics, Berlin, 1996
- [10] V.Girault,P.A.Raviart  
*Finite Element Approximation of the Navier-Stokes Equations*, Springer Verlag, Berlin, 1979
- [11] T.Grauschopf, M.Griebel, H.Regler  
*Additive Multilevel Preconditioners based on Bilinear Interpolation, Matrix Dependent Geometric Coarsening and Algebraic Multigrid Coarsening for Second Order Elliptic PDEs*, Technische Universität München, Institut für Informatik, SFB Bericht Nr.342/02/96 A, 1996
- [12] C.Großmann, H.-G.Roos  
*Numerik partieller Differentialgleichungen*, Teubner Studienbücher Mathematik, 1993

- [13] M.Griebel  
*Multilevelmethoden als Iterationsverfahren über Erzeugendensystemen*, B.G.Teubner, Stuttgart, 1994
- [14] W.Hackbusch  
*Iterative Löser großer schwachbesetzter Gleichungssysteme*, Teubner Studienbücher Mathematik, 1993
- [15] W.Hackbusch  
*Multigrid Methods and Application*, Springer Verlag, Berlin-Heidelberg-New York, 1985
- [16] W.Hackbusch  
*Theorie und Numerik elliptischer Differentialgleichungen*, B.G. Teubner, Stuttgart, 1986
- [17] G.Haase, U.Langer  
*Multigrid Methoden*, Johannes Kepler Universität Linz, Institut für Mathematik, 1998
- [18] R.Hiptmair  
*Multigrid Method for Maxwell's Equations*, Universität Augsburg, Institut für Mathematik, Report No. 374, 1997
- [19] R.Hiptmair  
*Nonconforming Vector Valued Finite Elements*, Universität Augsburg, Institut für Mathematik, Report No. 375, 1997
- [20] M.Jung, U.Langer  
*Finite-Elemente-Methode. Eine Einführung für Ingenieurstudenten*, Johannes Kepler Universität Linz, Institut für Mathematik, 1996
- [21] F.Kickinger  
*Algebraic Multigrid for Discrete Elliptic Second-Order Problems*, Johannes Kepler Universität Linz, Institut für Mathematik, Institutsbericht Nr.513, 1997
- [22] A.Kost  
*Numerische Methoden in der Berechnung elektromagnetischer Felder*, Springer-Verlag, Berlin, 1995
- [23] M.Kuhn  
*Efficient Parallel Numerical Simulations of Magnetic Field Problems* Dissertation, Johannes Kepler Universität Linz, Institut für Mathematik, 1998
- [24] U.Langer  
*Numerik I*, Skriptum zur gleichnamigen Vorlesung, Johannes Kepler Universität, Institut für Mathematik, Linz, 1996
- [25] U.Langer  
*Numerik II*, Skriptum zur gleichnamigen Vorlesung, Johannes Kepler Universität, Institut für Mathematik, Linz, 1996

- [26] P.Oswald  
*Multilevel Finite Element Approximation*, B.G.Teubner, Stuttgart, 1994
- [27] A.Reusken  
*On the Approximate Cyclic Reduction Preconditioner*, Rheinisch-Westfälische Technische Hochschule Aachen, Institut für Geometrie und Praktische Mathematik, Bericht Nr.144, 1997
- [28] J.W.Ruge, K.Stüben  
*Algebraic Multigrid (AMG)*, in Multigrid Methods (St.Mc Cormick ed.), Frontiers in Applied Mathematics, Vol.5, SIAM, Philadelphia, 1986
- [29] J.W.Ruge, K.Stüben  
*Efficient solution of finite difference and finite element equations*, in Multigrid Methods for integral and differential equations (D.Paddon and H.Holstein ed.), The Institute of Mathematics and it's Applications, Conference Series, New Series 3, Clarendon Press, Oxford, 1985
- [30] J.Schöberl  
*Robust Multigrid Preconditioning for Parameter-Dependent Problems I: The Stokes-type Case*, Johannes Kepler Universität Linz, Institut für Mathematik, Institutsbericht, 1997
- [31] J.Schöberl  
*FEPP - Finite Element ++*, [www.nathan.uni-linz.ac.at/Staff/joachim/cpp/doc/index.html](http://www.nathan.uni-linz.ac.at/Staff/joachim/cpp/doc/index.html), Johannes Kepler Universität Linz, Institut für Mathematik, 1997
- [32] O.Scherzer  
*Mathematische Modellierung*, Skriptum zur gleichnamigen Vorlesung, Johannes Kepler Universität, Ordinariat für Industriemathematik, Linz, 1997
- [33] P.Vanek, J.Mandel, M.Brezina  
*Algebraic Multigrid by Smoothed Aggregation for Second and Fourth Order Elliptic Problems*, Computing 56, 1996
- [34] W.L.Wendland (Editor)  
*Boundary Element Topics*, Preceedings of the Final Conference of the Priority Research Programme Boundary Element Methods 1989-1995 of the German Research Foundation, Springer Verlag, Berlin-Heidelberg-New York, 1997
- [35] P.M.de Zeeuw  
*Matrix Dependent Prolongations and Restrictions in a Black Box Multigrid*, J. Comp. and Appl. Mathematics 33, 1990
- [36] A.Zenisek, M.Vanmaele  
*The interpolation theorem for narrow quadrilateral isoparametric finite elements*, Numerische Mathematik 72, 1995