

Optimization of Photonic Band Structures

Zur Erlangung des akademischen Grades eines

DOKTORS DER NATURWISSENSCHAFTEN

von der Fakultät für Mathematik des
Karlsruher Instituts für Technologie (KIT)

genehmigte

DISSERTATION

von

Dipl.–Math. techn. Markus Richter

aus Frankfurt am Main

Tag der mündlichen Prüfung: 10. November 2010

Referent: Prof. Dr. Willy Dörfler

Koreferent: Prof. Dr. Christian Wieners

Preface

This dissertation is the result of a research project I conducted at the Institute for Applied Mathematics and Numerical Analysis of the Karlsruhe Institute of Technology (KIT) in the years 2006 to 2010. During that time I worked as a research assistant at the institute under the supervision of Prof. Dr. Willy Dörfler. The project was embedded within the research program of the DFG Research Training Group 1294 “Analysis, Simulation and Design of Nanotechnological Processes” at the Department of Mathematics.

At this point I would like to take the opportunity to thank a number of people, who supported me during my research. First and foremost I would like to thank my advisors, Prof. Dr. Willy Dörfler and Prof. Dr. Christian Wieners, for giving me the opportunity to work on an interesting and challenging problem, and for guiding me along my way towards this dissertation. I am particularly thankful for the many times I was able to discuss my ideas with them, or to seek their advice. I would also like to mention Prof. Wieners’ personal efforts in helping me with his numerical software library M++.

Furthermore, I would like to thank my colleagues at the institute and my fellow members of the Research Training Group 1294 for the good collaboration and for the many interesting discussions. I am especially indebted to Dr. Alexander Buloviyatov, who implemented some numerical algorithms, on which I could base my own ones. I also would like to mention Branimir Anić, Markus Bürg, Dr. Thomas Dohnal, Dr. Christian Engström, Markus Feist, Dr. Thomas Gauss, Dr. Vu Hoang, Florian Keller, and Daniel Maurer. A special thanks goes to Benjamin Exner, André Fuetterer, Dr. Desirée Hilbring, and Andreas Kuhn for proofreading parts of this dissertation.

I thank my family for their constant encouragement and support. I would also like to thank my closest friends for their company and for the many hours spent talking and laughing together. Finally, I thank Melanie for her kind, loving support.

Markus Richter
Karlsruhe, October 2010

Contents

1	Introduction	6
1.1	Motivation	6
1.2	Aims of This Work	9
1.3	Literature Review	10
1.4	Outline	12
2	Preliminaries	13
2.1	Notations and Conventions	13
2.2	The Cross Product	15
2.3	Local Lipschitz Continuity	16
2.4	Symmetry	17
2.5	Periodicity	19
3	The Mathematical Model	21
3.1	Modelling Photonic Crystals	21
3.2	Crystal Symmetries	23
3.3	Wave Propagation in Linear Dielectrics	30
3.4	Bloch Modes in Periodic Media	37
3.5	The Two-Dimensional Case	40
4	Spectral Theory	45
4.1	The Formal Setting	45
4.2	Sobolev Spaces of Periodic Functions	46
4.3	The Weak Formulation	59
4.4	Riesz–Schauder Theory	65
4.5	Auchmuty’s Principle	72
4.6	Photonic Band Structures and Band Diagrams	77
4.7	The Two-Dimensional Case	84
5	Photonic Band Structure Optimization	90
5.1	The Formal Setting	90

5.2	Regularity of the Goal Functionals	92
5.3	Existence of Optima in the TM Setting	96
5.4	Existence of Optima in Other Settings	103
5.5	Optimization Goals	107
6	Nonsmooth Analysis	110
6.1	Generalized Differentials	110
6.2	Generalized Differential Calculus	116
6.3	Differentiability of the Gap Width Functionals	120
6.4	The Two-Dimensional Case	125
7	A Generalized Gradient Method	127
7.1	Basic Concepts	127
7.2	Discretization	130
7.3	Incorporating Optimization Constraints	133
7.4	Preserving Symmetries	136
7.5	The Algorithm	138
7.6	Choosing Descent Directions	140
7.7	Remarks	147
8	A Level-Set Method	148
8.1	Basic Concepts	148
8.2	Motivation	150
8.3	The Algorithm	154
8.4	Remarks	157
9	Numerical Results	158
9.1	Maximizing TM Band Gaps	158
9.1.1	Results of the Generalized Gradient Method	159
9.1.2	Results of the Level Set Method	163
9.2	Maximizing TE Band Gaps	168
9.2.1	Results of the Generalized Gradient Method	168
9.2.2	Results of the Level Set Method	171
9.3	Maximizing Band Gaps of a 3D Photonic Crystal	173
10	Summary, Conclusions and Outlook	176
10.1	Summary	176
10.2	Comparison of the Optimization Algorithms	178
10.3	Open Problems	178
10.4	Final Remarks and Outlook	179

<i>CONTENTS</i>	5
A A FEM Toolbox for MATLAB	180
A.1 General Remarks	180
A.2 Data Structures and Algorithms	182
A.3 Some Customized Features	185
A.4 Implementation Examples	186
Frequently Used Symbols	191
About The Author	192
Bibliography	194

Chapter 1

Introduction

1.1 Motivation

Photonic crystals are materials, which are composed of two or more different dielectrics or metals, and which exhibit a spatially periodic structure, typically at the length scale of hundred nanometers. Depending on whether the periodicity extends into one, two or three space dimensions, a photonic crystal is called one-, two-, three-dimensional. Photonic crystals can be fabricated using nano-technological processes such as photolithography or vertical deposition methods. They also occur in nature, e.g. in the microscopic structure of certain bird feathers, butterfly wings, or beetle shells (see e.g. [12], [46]).

A characteristic feature of photonic crystals is that they strongly affect the propagation of light waves at certain optical frequencies. This is due to the fact that the optical density inside a photonic crystal varies periodically on the length scale of about 400 to 800 nanometers. One finds that the so-called *optical wavelengths* of light waves lie in precisely the same length scale. Light waves that penetrate a photonic crystal, are therefore subject to periodic, multiple diffraction, which leads to coherent wave interference inside the crystal. Depending on the frequency of the incident light wave this interference can either be constructive or destructive. In the latter case the light wave is not able to propagate inside the photonic crystal at all. Typically, this phenomenon only occurs for a bounded range of optical wave frequencies, if it does occur at all. Such a range of inhibited wave frequencies is called a *photonic band gap*. Light waves with frequencies inside a photonic band gap are totally reflected by the photonic crystal. It is this effect, which causes e.g. to the iridescent colors of peacock feathers (see [54]).

Whether or not photonic band gaps occur, strongly depends on the geometric structure of the photonic crystal, as well as on the contrast in optical density between the different materials the photonic crystal is built of. Apart from photonic

band gaps a photonic crystal may also exhibit other optical phenomena, such as a large refractive index or the ability to slow down the group velocity of light pulses considerably. These phenomena are also caused by periodic diffraction and coherent wave interference inside the crystal.

In the late 1980s Eli Yablonovitch and Sajeev John discovered that the band gap of a photonic crystal can be used to inhibit spontaneous photon emissions (see [43], [80]). It was also then that the term *photonic crystal* was coined. About a century earlier Lord Rayleigh studied the electromagnetic transmissibility of multi-layered dielectrics, and found that for certain frequencies of the incident wave total reflection can occur (see [63]). This phenomenon, which also stems from coherent wave interference inside the multi-layered material, was used to create frequency-selective optical mirrors, the so-called *distributed Bragg reflectors*. In essence these devices consist of periodically alternating layers of different dielectric materials. From a modern perspective, distributed Bragg reflectors can be viewed as one-dimensional photonic crystals. Two- and three-dimensional photonic crystals can be used to guide light waves along a defect in the periodic structure. It is expected that photonic crystals will play a key role in the development of nanometer-scale all-optical communication devices, such as nanometer wave-guides, optical multiplexers, logical gates or frequency filters. Some of these devices could already been realized in the laboratory. However, an industry-scale production is not yet in sight.

Many optical properties of a photonic crystal are determined by its so-called *photonic band structure*. The band structure consists of countably many of so-called *photonic bands*. Each band discloses the dispersion relation of a time-harmonic electromagnetic wave, which is able to propagate in the crystal. Mathematically a photonic band is the graph of a function $\mathbf{k} \mapsto \omega(\mathbf{k})$, where the vector \mathbf{k} varies over a compact subset \mathbb{B} of \mathbb{R}^3 , the so-called first *Brillouin zone*. For every vector $\mathbf{k} \in \mathbb{B}$ the corresponding function value $\omega(\mathbf{k})$ is given by a solution of an eigenvalue problem, which is of the form

$$\begin{cases} (\nabla + i\mathbf{k}) \times [\rho(\nabla + i\mathbf{k}) \times \mathbf{u}] = \omega^2 \mathbf{u} & \text{in } \Omega, \\ (\nabla + i\mathbf{k}) \cdot \mathbf{u} = 0 & \text{in } \Omega. \end{cases} \quad (1.1)$$

Here, Ω denotes a *fundamental cell of periodicity* of the photonic crystal. A fundamental cell of periodicity is a bounded section of space within a photonic crystal. The defining property of a fundamental cell is that the entire photonic crystal structure can be reproduced by repeating the structure within the fundamental cell periodically. The crystal structure itself is represented by the real-valued function ρ . Typically ρ is discontinuous and takes only a finite number of different function values. The eigenfunctions \mathbf{u} are required to satisfy periodic boundary conditions. It can be shown, that for every vector $\mathbf{k} \in \mathbb{B}$ there exists a countable

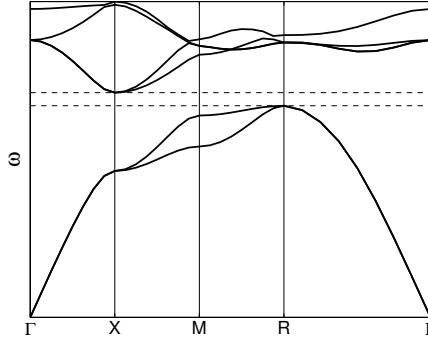


Figure 1.1: Photonic band structure exhibiting a band gap (dashed lines) between the second and third photonic band.

number of positive, real numbers $\omega_1(\mathbf{k}) \leq \omega_2(\mathbf{k}) \leq \dots$, etc., such that $\omega_j(\mathbf{k})^2$ is an eigenvalue of (1.1) for every $j \in \mathbb{N}$. Given an index $j \in \mathbb{N}$, the graph of the function $\mathbf{k} \mapsto \omega_j(\mathbf{k})$ is commonly referred to as the j -th *photonic band*, a photonic band gap exists between the j -th and the $(j + 1)$ -th photonic band, if and only if

$$\min_{\mathbf{k} \in \mathbb{B}} \omega_{j+1}(\mathbf{k}) - \max_{\mathbf{k} \in \mathbb{B}} \omega_j(\mathbf{k}) > 0. \quad (1.2)$$

The so-called *gap width* is given by the left-hand side of (1.2). Photonic band structures are usually depicted in so-called *photonic band diagrams*. In such a band diagram each photonic band is represented by a curve. Photonic band gaps are revealed as horizontal regions in the diagram, which are not crossed by any curve (see Figure 1.1). The vertical extent of such a region corresponds to the gap width.

For many applications, it is desirable to have photonic crystals with large band gaps as large as possible. Given two different materials, the question arises how a photonic crystal should be built from these two materials in order to exhibit a maximal band gap. Recall that the structure of a photonic crystal is represented by the function ρ in the eigenvalue problem (1.1). A photonic crystal, which consists of exactly two different materials, is represented by a two-valued function ρ on Ω with values ρ_0 and ρ_1 . Conversely, every two-valued function on Ω with values ρ_0 and ρ_1 represents the structure of a photonic crystal built from the same two materials. We shall call such functions *admissible functions*. For every admissible function ρ and every vector $\mathbf{k} \in \mathbb{B}$, one can solve (1.1) and obtain solutions $\omega_1(\rho, \mathbf{k}) \leq \omega_2(\rho, \mathbf{k}) \leq \dots$, etc., which determine the band structure of the photonic crystal represented by ρ . Given a fixed index $j \in \mathbb{N}$ the width of the band gap between the j -th and $(j + 1)$ -th band of that crystal, in case it exists, is

given by

$$w_j(\rho) := \min_{\mathbf{k} \in \mathbb{B}} \omega_{j+1}(\rho, \mathbf{k}) - \max_{\mathbf{k} \in \mathbb{B}} \omega_j(\rho, \mathbf{k}).$$

With this, the problem of finding a photonic crystal structure, which exhibits a maximal band gap between the j -th and $(j+1)$ -th photonic band can be formulated as follows: Find an admissible function ρ^* , such that $w_j(\rho^*) \geq w_j(\rho)$ for every admissible function ρ . A problem of this type is commonly referred to as a *photonic band gap maximization problem* (hereafter abbreviated as PBGM).

More generally, let J be a functional, such that for every admissible function ρ the corresponding value $J(\rho)$ of the functional is determined by the functions $\mathbf{k} \mapsto \omega_1(\rho, \mathbf{k})$, $\mathbf{k} \mapsto \omega_2(\rho, \mathbf{k})$, \dots , etc. An optimization problem could then consist in finding an admissible function ρ^* , such that $J(\rho^*) \leq J(\rho)$ for every admissible function ρ . We shall call a problem of this type a *photonic band structure optimization problem* (hereafter abbreviated as PBSOP). The functional J is commonly referred to as the *goal functional* of the problem. Clearly, photonic band gap maximization problems are those PBSOPs, whose goal functionals are given by $J = -w_j$ for some $j \in \mathbb{N}$.

Photonic band structure optimization problems present a number of interesting challenges for a mathematician. In many relevant cases the goal functional J fails to be any more regular than Lipschitz continuous. Moreover, the set of admissible functions is infinite-dimensional and non-compact, which makes it difficult to prove the existence of optimal solutions. Computing the band structure of a photonic crystal is a challenging task from a numerical point of view, especially when the crystal is three-dimensional. One of the reasons for this is that admissible functions, which are by definition two-valued, enter as discontinuous coefficients in the eigenvalue problem (1.1). Therefore, sophisticated numerical schemes need to be employed in order to obtain accurate results. Finally, the development and analysis of numerical optimization algorithms for PBSOPs provides a field of research that mathematicians can contribute to.

In this work we discuss in detail the mathematical theory, which relates the geometric structure of a photonic crystal to its band structure, and develop different optimization algorithms for PBSOPs. A more general introduction to photonic crystals can be found in the books [42], [60], or [66]. A survey on mathematical concepts related to photonic crystals is given in [50].

1.2 Aims of This Work

The aim of this thesis is the mathematical study of photonic band structure optimization problems (PBSOPs), as well as the development of suitable optimization algorithms for these problems. Moreover, we aim at finding solutions of photonic

band gap maximization problems (PBGMPs) for two- and three-dimensional photonic crystals.

1.3 Literature Review

Photonic band structure optimization problems are, in essence, eigenfrequency optimization problems. Problems of this type also occur e.g. in structural engineering and acoustics. As an example, we mention the problem of finding the optimal design of a two-density composite membrane, where the optimization goal typically is the minimization or maximization of certain eigenfrequencies. This problem was studied in detail by Cox and McLaughlin (cf. [31], [32]), who proved the existence of optimal designs and characterized them in terms of certain level-sets of eigenfunctions. Some fifty years earlier, Krein was able to identify two-density composite strings with minimal and maximal eigenfrequencies (cf. [48]). The eigenvalue problem in both cases is of the form

$$\begin{cases} -\Delta u = \lambda \rho u & \text{in } \Omega, \\ u = 0 & \text{on } \partial\Omega, \end{cases}$$

where Ω is a bounded domain. The function ρ is two-valued on Ω and models the density distribution of the string or membrane. As will turn out in Section 3.5, eigenvalue problems of similar type also arise for certain PBSOPs involving one- and two-dimensional photonic crystals. In Section 5.4 we will show how some results found by Cox, McLaughlin and Krein carry over to these problems.

Another optimization problem, which is in a way related to PBSOPs, is an optimal design problem concerning the temperature distribution in two-phase conductors. Under certain model assumptions the temperature distribution is determined by an eigenvalue problem of the form

$$\begin{cases} -\nabla \cdot [\rho \nabla u] = \lambda u & \text{in } \Omega, \\ u = 0 & \text{on } \partial\Omega, \end{cases}$$

where ρ , again, is a two-valued function on a bounded domain Ω . The problem was analysed by Cox and Lipman (cf. [27]) and further studied by Alvio, Lion and Trombetti (cf. [5]), as well as by Conca, Mahadevan and Sanz (cf. [26]). Obviously, the eigenvalue problem (1.3) bears some similarity to the the eigenvalue problem (1.1), which determines the band structure of a photonic crystal. In Section 5.4 we will discuss the findings of the authors mentioned above and illustrate how these findings relate to PBSOPs.

Photonic band gap maximization problems have already been considered in the literature, though only for two-dimensional photonic crystals. Cox and Dobson developed an optimization algorithm based on a generalized gradient descent method

(cf. [29]). This algorithm served as a prototype for the optimization algorithms we develop in Chapter 7. It should be noted, however, that the optimization problems considered by Cox and Dobson differ slightly from PBGMPs as defined in the previous Section 1.1. In their paper the authors choose a fixed frequency lying inside a photonic band gap and aim at pushing the nearest photonic bands away from that frequency. In doing so Cox and Dobson are able to find photonic crystals, which exhibit a large band gap with a prescribed frequency at the center. In a second paper Cox and Dobson refine their algorithm allowing them to open up photonic band gaps between specific bands (cf. [30]).

Using a level-set method (see e.g. [58]), Kao, Osher and Yablonovitch were able to find two-dimensional photonic crystals with maximal band gaps (cf. [44]). While providing an extensive list of optimized crystal structures, the authors do not establish a rigorous analysis of their method. A non-rigorous justification can be found, however, in a paper by Osher and Santosa on the optimization of two-density composite membranes via level-set methods (cf. [59]).

Yet another approach was presented by Sigmund and Søndergaard, who treated PBGMPs and PBSOPs as topology optimization problems (see e.g. [11]). In a paper on the systematic design of two-dimensional photonic crystals (cf. [69]) the authors compute optimized crystal structures using the method of moving asymptotes (see [71]). The authors remark, however, that this method could only be applied to rather coarse discretizations of the problem due to computing-time and storage limitations. Preble, M. Lipson, and H. Lipson used evolutionary algorithms in order to design two-dimensional photonic crystals with maximal band gaps (cf. [61]). Each evolutionary algorithm needed to create more than 1000 generations in order to find optimized crystal structures, which is why the primitive cell of each photonic crystal was discretized by a 32×32 grid only.

As we mentioned before, there are no published results neither on PBGMPs nor on PBSOPs for three-dimensional photonic crystals so far. A possible reason for this may be the fact that these problems are computationally more demanding than their two-dimensional analogues. As will be discussed in Section 3.5, photonic band structures of two-dimensional crystals are determined by Laplace-type or second-order divergence-type eigenvalue problems. Such eigenvalue problems can be discretized using standard finite difference or finite element methods, and band structure computations can be carried out on an average computer. The situation is much different for three-dimensional photonic crystals, whose band structures are determined by second-order curl-curl-type eigenvalue problems of the form (1.1). In order to obtain suitable discretizations of these problems, specialized finite element schemes have to be employed. Moreover, efficient numerical algorithms and high-performance computer architectures are generally needed in order to carry out band structure computations within reasonable time. Thus, the

numerical results presented in Chapter 9 were only made possible by the efforts of Bulovyatov (cf. [17]), who developed a parallel multi-grid algorithm for band structure computations of three-dimensional photonic crystals based on a finite element library developed by Wieners (see [76], [77]).

1.4 Outline

This work is organized as follows. In Chapter 2 we present the notation, certain function spaces and fundamental concepts that are used in this work. In Chapter 3 we introduce the mathematical model for wave propagation in photonic crystals, focussing on the underlying physical principles and model assumptions. A family of eigenvalue problems governing the propagation of so-called Bloch modes in photonic crystals is given. The formulation of the eigenvalue problems is made mathematically precise in Chapter 4. There, we also provide the appropriate spectral theory and introduce the band structure of a photonic crystal, as well as photonic band gaps. The chapter is completed by a discussion of the Floquet–Bloch theory, which relates the band structure of a photonic crystal to its transmissibility for electromagnetic radiation. In Chapter 5 we address the problem of optimizing the geometrical structure of a photonic crystal with respect to certain properties of its band structure. We establish a generic minimization problem and analyse its properties. In Chapter 5 we also state some model problems, which motivated this work. In Chapter 6 we review important concepts of nonsmooth analysis. In particular, the concept of generalized differentials is introduced. In Chapter 7 we develop an algorithm for the numerical solution of photonic band structure optimization problems. In essence, this algorithm is based on the idea of gradient descent methods. A slight drawback inherent to the algorithm is that it fails to produce admissible crystal structures in general. In Chapter 8 we therefore develop another optimization algorithm based on a level-set method. In Chapter 9 we present and discuss some numerical results of the algorithms. Conclusions and an outlook on future work are given in Chapter 10.

Chapter 2

Preliminaries

In this chapter we introduce some fundamental concepts, notations and conventions that will be used throughout the rest of this work. Basic notations and conventions are introduced in Section 2.1. Important results concerning the vector cross product are listed in Section 2.2. In Section 2.3 we discuss the notion of local Lipschitz continuity. In Section 2.4 we review the concept of symmetry and related aspects of group theory. Periodicity is a special case of symmetry and will be discussed in Section 2.5. In both Sections 2.4 and 2.5 we introduce fundamental mathematical structures, such as symmetry groups or Bravais lattices, which will be used in Chapter 3 to build a mathematical model for the medium structures of photonic crystals.

2.1 Notations and Conventions

In this work we denote by \mathbb{N} the set of natural numbers not including zero, and by \mathbb{N}_0 the set $\mathbb{N} \cup \{0\}$. Euler's number and the imaginary unit are denoted by the upright letters e and i , respectively. We use boldface letters to denote vector- or matrix-valued quantities. Unless stated otherwise, \mathbf{u}_i denotes the i -th component of the vector-valued quantity \mathbf{u} for a given index $i \in \mathbb{N}$. Similarly, \mathbf{A}_{ij} denotes the component in the i -th row and j -th column of the matrix-valued quantity \mathbf{A} for given indices $i \in \mathbb{N}$ and $j \in \mathbb{N}$.

Given a complex-valued quantity z , we denote by \bar{z} its complex conjugate. Sesquilinear forms on a complex vector space are understood to be conjugate-linear in the first (left-hand) argument and linear in the second (right-hand) argument. This applies, in particular, to the standard inner product on \mathbb{C}^n , $n \in \mathbb{N}$, which we denote by $\langle \cdot, \cdot \rangle$. In contrast to this, we denote by \cdot the dot product on \mathbb{C}^n . The connection between the inner product and the dot product is given by $\langle \mathbf{x}, \mathbf{y} \rangle = \bar{\mathbf{x}} \cdot \mathbf{y}$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$. If X is a real or complex topological vector space,

then X^* denotes its topological dual, i.e., the space X^* consists of all continuous linear functionals on X . Given another topological vector space Y and an operator $T : X \rightarrow Y$, we sometimes denote the image of a point $x \in X$ under T by Tx , if T is linear or conjugate-linear, or by $T[x]$, if T is linear. If X and Y are normed spaces with norms $\|\cdot\|_X$ and $\|\cdot\|_Y$, and if T is a continuous, linear or conjugate-linear operator, we denote by $\|T\|_{X,Y}$ the operator norm of T , given by

$$\|T\|_{X,Y} := \sup_{x \in X \setminus \{0\}} \frac{\|Tx\|_Y}{\|x\|_X}.$$

Given a set $S \subset X$, we denote by $\text{conv}(S)$ the open, convex hull of S and by $\overline{\text{conv}}(S)$ its closure.

Given a set $\Omega \subseteq \mathbb{R}^n$, where $n \in \mathbb{N}$, a Banach space Y with norm $\|\cdot\|_Y$, and a real number $p \in \mathbb{R}$, with $1 \leq p < \infty$, we define the norms $\|\cdot\|_{\Omega,p}$ and $\|\cdot\|_{\Omega,\infty}$ by

$$\begin{aligned} \|f\|_{\Omega,p} &:= \left(\int_{\Omega} \|f\|_Y^p \right)^{1/p} && \text{for all } f \in L^p(\Omega, Y), \\ \|f\|_{\Omega,\infty} &:= \text{ess sup}_{\Omega} (\|f\|_Y) && \text{for all } f \in L^\infty(\Omega, Y). \end{aligned}$$

As a convention, we further constitute that $\|\cdot\|_{\Omega} := \|\cdot\|_{\Omega,2}$. If Y is a Hilbert space with inner product $\langle \cdot, \cdot \rangle_Y$, we define the inner product $\langle \cdot, \cdot \rangle_{\Omega}$ by

$$\langle f, g \rangle_{\Omega} := \int_{\Omega} \langle f, g \rangle_Y \quad \text{for all } f, g \in L^2(\Omega, Y).$$

Whenever $F(\Omega, Y)$ denotes a set of functions mapping Ω into Y , we denote by $F(\Omega)$ the respective set of functions mapping Ω into \mathbb{C} . This convention applies, in particular, to the standard Lebesgue, Sobolev, and Hölder spaces $L^p(\Omega)$, $H^s(\Omega)$, and $C^{m,\alpha}(\Omega)$.

By $\mathbf{H}(\text{curl}; \Omega)$ and $\mathbf{H}(\text{div}; \Omega)$ we denote the Sobolev spaces, which consist of all functions from Ω into \mathbb{C}^3 that admit weak curl or divergence fields in $L^2(\Omega)^3$ or $L^2(\Omega)$, respectively. A detailed characterization of the Sobolev space $\mathbf{H}(\text{curl}; \Omega)$ can be found in Chapter 7, Section 4.1 in [34]. The Sobolev space $\mathbf{H}(\text{div}; \Omega)$ is discussed in detail in Section 2.1 in [38]. The spaces $\mathbf{H}(\text{curl}; \Omega)$, $\mathbf{H}(\text{div}; \Omega)$ and $H^1(\Omega)$ are equipped with the inner products $\langle \cdot, \cdot \rangle_{\text{curl}, \Omega}$, $\langle \cdot, \cdot \rangle_{\text{div}, \Omega}$, and $\langle \cdot, \cdot \rangle_{1, \Omega}$, respectively, which are given by

$$\begin{aligned} \langle \mathbf{u}, \mathbf{v} \rangle_{\text{curl}, \Omega} &:= \langle \nabla \times \mathbf{u}, \nabla \times \mathbf{v} \rangle_{\Omega} + \langle \mathbf{u}, \mathbf{v} \rangle_{\Omega} && \text{for all } \mathbf{u}, \mathbf{v} \in \mathbf{H}(\text{curl}; \Omega), \\ \langle \mathbf{u}, \mathbf{v} \rangle_{\text{div}, \Omega} &:= \langle \nabla \cdot \mathbf{u}, \nabla \cdot \mathbf{v} \rangle_{\Omega} + \langle \mathbf{u}, \mathbf{v} \rangle_{\Omega} && \text{for all } \mathbf{u}, \mathbf{v} \in \mathbf{H}(\text{div}; \Omega), \\ \langle p, q \rangle_{1, \Omega} &:= \langle \nabla p, \nabla q \rangle_{\Omega} + \langle p, q \rangle_{\Omega} && \text{for all } p, q \in H^1(\Omega). \end{aligned}$$

Here as in the following, we denote by $\nabla \times$, $\nabla \cdot$, and ∇ the curl, divergence and gradient operator. The norms, which are induced by the inner products, are denoted by $\|\cdot\|_{\text{curl}, \Omega}$, $\|\cdot\|_{\text{div}, \Omega}$, and $\|\cdot\|_{1, \Omega}$, respectively.

2.2 The Cross Product

In this work we quite perform vector calculations involving the cross product. Therefore, we find it useful to mention some important results related to this vector operation here.

First, let us recall some vector identities. An important identity, which concerns so-called *vector triple products*, reads

$$\mathbf{x} \times (\mathbf{y} \times \mathbf{z}) = (\mathbf{x} \cdot \mathbf{z})\mathbf{y} - (\mathbf{x} \cdot \mathbf{y})\mathbf{z} \quad \text{for all } \mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{C}^3. \quad (2.1)$$

For so-called *scalar triple products* we have the well-known identity

$$\mathbf{x} \cdot (\mathbf{y} \times \mathbf{z}) = \mathbf{y} \cdot (\mathbf{z} \times \mathbf{x}) = \mathbf{z} \cdot (\mathbf{x} \times \mathbf{y}) \quad \text{for all } \mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{C}^3, \quad (2.2)$$

which is often verbalized by saying the the scalar triple product is invariant under cyclic permutations of the vectors. Furthermore, *Lagrange's identity* states that

$$\langle \mathbf{x} \times \mathbf{y}, \tilde{\mathbf{x}} \times \tilde{\mathbf{y}} \rangle = \langle \mathbf{x}, \tilde{\mathbf{x}} \rangle \langle \mathbf{y}, \tilde{\mathbf{y}} \rangle + \langle \mathbf{y}, \tilde{\mathbf{x}} \rangle \langle \mathbf{x}, \tilde{\mathbf{y}} \rangle \quad \text{for all } \mathbf{x}, \tilde{\mathbf{x}}, \mathbf{y}, \tilde{\mathbf{y}} \in \mathbb{C}^3. \quad (2.3)$$

Letting $\tilde{\mathbf{x}} = \mathbf{x}$ and $\tilde{\mathbf{y}} = \mathbf{y}$ in (2.3), we find that

$$|\mathbf{x}|^2 |\mathbf{y}|^2 = |\langle \mathbf{x}, \mathbf{y} \rangle|^2 + |\mathbf{x} \times \mathbf{y}|^2 \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathbb{C}^3, \quad (2.4)$$

where $|\cdot|$ denotes the Euclidean norm on \mathbb{C}^3 . This norm identity in particular implies the following Cauchy–Schwarz-like inequality, which reads

$$|\mathbf{x} \times \mathbf{y}| \leq |\mathbf{x}| |\mathbf{y}| \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathbb{C}^3.$$

Next, we introduce a useful concept for representing the action of cross products on \mathbb{C}^3 . Given a vector $\mathbf{z} \in \mathbb{C}^3$, we define the matrix $[\mathbf{z}]_{\times} \in \mathbb{C}^{3 \times 3}$ by

$$[\mathbf{z}]_{\times} := \begin{pmatrix} 0 & -z_3 & z_2 \\ z_3 & 0 & -z_1 \\ -z_2 & z_1 & 0 \end{pmatrix}. \quad (2.5)$$

We refer to the matrix $[\mathbf{z}]_{\times}$ as the the *cross product matrix* of \mathbf{z} . One easily verifies that

$$[\mathbf{z}]_{\times} \mathbf{x} = \mathbf{z} \times \mathbf{x} \quad \text{for all } \mathbf{x} \in \mathbb{C}^3.$$

For $\mathbf{z} \in \mathbb{C}^3 \setminus \{\mathbf{0}\}$ one finds that the null space and the image space of $[\mathbf{z}]_{\times}$ are given by

$$\begin{aligned} \ker([\mathbf{z}]_{\times}) &= \text{span}\{\mathbf{z}\}, \\ \text{im}([\mathbf{z}]_{\times}) &= (\text{span}\{\mathbf{z}\})^{\perp}, \end{aligned}$$

where $(\text{span}\{\mathbf{z}\})^\perp$ denotes the orthogonal complement of $\text{span}\{\mathbf{z}\}$ in \mathbb{C}^3 .

Finally, we mention a vector identity concerning the behaviour of cross products under linear coordinate transformations. Given a matrix $\mathbf{A} \in \mathbb{C}^{3 \times 3}$, one can show by direct computation that

$$(\mathbf{A}\mathbf{x}) \times (\mathbf{A}\mathbf{y}) = \text{cof}(\mathbf{A})(\mathbf{x} \times \mathbf{y}) \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathbb{C}^3, \quad (2.6)$$

where $\text{cof}(\mathbf{A})$ denotes the *cofactor matrix* of \mathbf{A} . The cofactor matrix is given by

$$\text{cof}(\mathbf{A})_{ij} := (-1)^{i+j} \det(\mathbf{A}^{[i,j]}) \quad \text{for all } i, j = 1, 2, 3, \quad (2.7)$$

where $\mathbf{A}^{[i,j]}$ denotes the matrix, which is obtained by omitting the i -th row and the j -th column in the matrix \mathbf{A} for $i, j = 1, 2, 3$. Recall that the cofactor matrix $\text{cof}(\mathbf{A})$ is the transposed of the so-called *adjugate matrix* $\text{adj}(\mathbf{A})$ of \mathbf{A} . We thus have

$$\text{cof}(\mathbf{A}) = \det(\mathbf{A})\mathbf{A}^{-\text{T}} \quad \text{for all } \mathbf{A} \in \text{GL}_3(\mathbb{C}), \quad (2.8)$$

where $\text{GL}_3(\mathbb{C})$ denotes the set consisting of all invertible 3×3 matrices with complex components.

2.3 Local Lipschitz Continuity

In this section we review the concept of locally Lipschitz continuous functions on Banach spaces. It will turn out that such functions play an important role in optimization problems involving band structures of photonic crystals.

Given two real or complex Banach spaces X and Y with norms $\|\cdot\|_X$ and $\|\cdot\|_Y$, a function $f : X \rightarrow Y$ is called *Lipschitz continuous*, if there exists a non-negative number $c_f \geq 0$, which depends on f only, such that

$$\|f(x) - f(y)\|_Y \leq c_f \|x - y\|_X \quad \text{for all } x, y \in X. \quad (2.9)$$

Given a constant $c > 0$, we shall call the function f above *c-Lipschitz continuous*, if (2.9) holds for $c_f = c$. Every number $c_f \geq 0$, which satisfies (2.9), is called a *Lipschitz constant* of f . It is well-known that Lipschitz continuity implies uniform continuity.

Given a point $x_0 \in X$, a function $f : X \rightarrow Y$ is called *Lipschitz continuous near x_0* , if there exists a neighbourhood $U \subseteq X$ of x_0 , such that f is Lipschitz continuous on U , i.e., if there exists a constant $c_{f,U} > 0$, which only depends on f and the neighbourhood U , such that

$$\|f(x) - f(y)\|_Y \leq c_{f,U} \|x - y\|_X \quad \text{for all } x, y \in U. \quad (2.10)$$

Given a constant $c > 0$, we shall call the function f above *c-Lipschitz continuous near x_0* , if there exists a neighbourhood U of x_0 , such that (2.10) holds for $c_{f,U} = c$. A function $f : X \rightarrow Y$ is called *locally Lipschitz continuous*, if it is Lipschitz continuous near every point $x \in X$.

Clearly, every Lipschitz continuous function is also locally Lipschitz continuous. It is also easy to see that local Lipschitz continuity implies continuity. It should be noted, however, that local Lipschitz continuity does not imply uniform continuity. To see this, consider e.g. the function $x \mapsto 1/x$ from $\mathbb{R}_{>0}$ into \mathbb{R} . Using the Mean-Value Theorem, one can show that this function is locally Lipschitz continuous, but not uniformly continuous. Finally, we remark that a function defined on a compact set is Lipschitz continuous if and only if it is locally Lipschitz continuous.

2.4 Symmetry

In this section we briefly review some mathematical concepts related to symmetry. As is well-known, symmetry can be described in terms of groups of isometries acting on metric spaces. Recall that a group G with neutral element e acts on a non-empty set X by virtue of a mapping $G \times X \rightarrow X$, which maps each pair $(g, x) \in G \times X$ to an element $x^g \in X$, and which satisfies $x^{gh} = (x^g)^h$, as well as $x^e = x$ for all $g, h \in G$ and all $x \in X$. This mapping $(g, x) \mapsto x^g$ is called the *group action of G on X* . Given a group element $g \in G$, an element $x \in X$ is said to be *invariant under the action of g on X* , if $x^g = x$. Furthermore, x is said to be *invariant under the action of G on X* , if $x^g = x$ for all $g \in G$. For each element $x \in X$ the set $G_x := \{g \in G \mid x^g = x\}$ is called the *stabilizer of x in G* . The stabilizer of x in G is a subgroup of G (notation: $G_x \leq G$). More precisely, G_x is the maximal subgroup of G under whose action on X the element x is invariant. In general, G_x is not a normal subgroup of G .

Given a metric space X , we denote by $\text{Iso}(X)$ the *isometry group of X* , which consists of all isometries from the metric space X onto itself with the function composition as group operation. The isometry group $\text{Iso}(X)$ canonically acts on X , on the power set $\mathcal{P}(X)$, and on every set of functions from X into another set Y . The respective group actions $\text{Iso}(X) \times X \rightarrow X$, $\text{Iso}(X) \times \mathcal{P}(X) \rightarrow \mathcal{P}(X)$, and $\text{Iso}(X) \times Y^X \rightarrow Y^X$ are given by

$$\begin{aligned} x^\varphi &:= \varphi(x) && \text{for all } \varphi \in \text{Iso}(X), x \in X, \\ S^\varphi &:= \varphi(S) && \text{for all } \varphi \in \text{Iso}(X), S \in \mathcal{P}(X), \\ f^\varphi &:= f \circ \varphi && \text{for all } \varphi \in \text{Iso}(X), f \in Y^X. \end{aligned}$$

It is well-known that the isometry group $\text{Iso}(\mathbb{R}^n)$ of the n -dimensional Euclidean space, where $n \in \mathbb{N}$, consists exactly of those affine mappings $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^n$, which

are of the form

$$\varphi(\mathbf{x}) := \Theta \mathbf{x} + \mathbf{a},$$

where $\Theta \in O_n(\mathbb{R})$ is an orthogonal $n \times n$ matrix with components in \mathbb{R} , and where $\mathbf{a} \in \mathbb{R}^n$ is a vector. Depending on the nature of Θ and \mathbf{a} , specific elements φ of $\text{Iso}(\mathbb{R}^n)$ can be identified e.g. as *rotations*, *reflexions*, *inversions*, or *translations* in \mathbb{R}^n .

Given a subgroup G of the isometry group $\text{Iso}(X)$ of an arbitrary metric space X , a subset S of X is called *symmetric with respect to G* or simply *G -symmetric*, if S is invariant under the action of G on $\mathcal{P}(X)$, i.e., if

$$\varphi(S) = S \quad \text{for all } \varphi \in G.$$

Similarly, a function $f \in Y^X$, where Y is a non-empty set, is called *G -symmetric*, if f is invariant under the action of G on Y^X , i.e., if

$$f \circ \varphi = f \quad \text{for all } \varphi \in G.$$

The corresponding stabilizers $\text{Iso}(X)_S$ and $\text{Iso}(X)_f$ are called the *symmetry groups* of the set S and the function f , respectively. Elements which belong to the symmetry group of an object are called *symmetry operations* of the object. If $X = \mathbb{R}^n$ for some $n \in \mathbb{N}$, one can think of a symmetry operation as a Euclidean motion, which takes an object into itself. The object's symmetry group is precisely the subgroup of $\text{Iso}(\mathbb{R}^n)$, which consists of all such operations.

Given a subgroup G of $\text{Iso}(X)$, a connected subset F of X is called a *fundamental region* of G , if it satisfies

$$\begin{aligned} \varphi(F) \cap \psi(F) &= \emptyset && \text{for all } \varphi, \psi \in G \text{ with } \varphi \neq \psi, \\ \bigcup_{\varphi \in G} \varphi(F) &= X. \end{aligned}$$

Hence, a fundamental region of G is a set whose images under the actions of G cover the entire metric space without overlapping. It follows that every G -symmetric function is determined by its restriction to a fundamental region of G . The interior of a fundamental region of G is called a *fundamental domain* of G .

Given a finite subgroup G of $\text{Iso}(X)$ and a function $f \in Y^X$, where Y is a vector space over \mathbb{Q} , we define the function $f^{(G)} \in Y^X$ by

$$f^{(G)} := \frac{1}{|G|} \sum_{\varphi \in G} (f \circ \varphi). \quad (2.11)$$

The function $f^{(G)}$ is called the *G -symmetrization of f* . Obviously, for every function $f \in Y^X$ the corresponding G -symmetrization $f^{(G)}$ is a G -symmetric function. Furthermore, one can show that $f^{(G)}$ coincides with f , if and only if f is G -symmetric.

2.5 Periodicity

Periodicity is a special case of symmetry, namely that of symmetry with respect to a finitely generated group of translations. Given a vector $\mathbf{a} \in \mathbb{R}^n$, $n \in \mathbb{N}$, the mapping $\tau_{\mathbf{a}} : \mathbb{R}^n \rightarrow \mathbb{R}^n$, which is defined by

$$\tau_{\mathbf{a}}(\mathbf{x}) := \mathbf{x} + \mathbf{a} \quad \text{for all } \mathbf{x} \in \mathbb{R}^n \quad (2.12)$$

is called the *translation in \mathbb{R}^n by \mathbf{a}* . One easily verifies that the set of all translations in \mathbb{R}^n constitutes a normal subgroup of $\text{Iso}(\mathbb{R}^n)$, the isometry group of \mathbb{R}^n . A discrete subset Λ of \mathbb{R}^n , which is of the form

$$\Lambda = \{z_1 \mathbf{a}^{(1)} + \cdots + z_r \mathbf{a}^{(r)} \mid z_1, \dots, z_r \in \mathbb{Z}\},$$

where $\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(r)} \in \mathbb{R}^n$ are linearly independent vectors, is called a *Bravais lattice* or simply a *lattice in \mathbb{R}^n of rank r* . The vectors $\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(r)}$, which form a \mathbb{Z} -linear basis of Λ , are often referred to as *primitive vectors of Λ* . Here as in the following, the term \mathbb{Z} -linear refers to the algebraic structure of a vector space over the set of integer numbers \mathbb{Z} . The elements of a lattice are called *lattice points* or *lattice vectors*. Note that every lattice contains the origin of the coordinate system as a lattice point.

If Λ_1 and Λ_2 are two Bravais lattices in \mathbb{R}^n , such that $\Lambda_1 \subseteq \Lambda_2$, we call Λ_1 a *sublattice of Λ_2* and write $\Lambda_1 \leq \Lambda_2$. Clearly, Λ_1 is a sublattice of Λ_2 , if and only if there exists a set of primitive vectors of Λ_1 , such that every vector is equal to a \mathbb{Z} -linear combination of primitive vectors of Λ_2 .

Every Bravais lattice $\Lambda \subset \mathbb{R}^n$ is a finitely generated, free, Abelian additive subgroup of \mathbb{R}^n . The corresponding set of *lattice translations*

$$\text{Trn}(\Lambda) := \{\tau_{\mathbf{a}} \mid \mathbf{a} \in \Lambda\} \quad (2.13)$$

is a finitely generated, free, Abelian subgroup of the isometry group $\text{Iso}(\mathbb{R}^n)$. Clearly, every sublattice of Λ is a normal subgroup of Λ , and the corresponding subgroup of lattice translation is normal in $\text{Trn}(\Lambda)$. Furthermore, we remark that every lattice $\Lambda \subset \mathbb{R}^n$ defines an equivalence relation \equiv_{Λ} on \mathbb{R}^n , which is given by

$$\mathbf{x} \equiv_{\Lambda} \mathbf{y} : \iff \mathbf{x} - \mathbf{y} \in \Lambda \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathbb{R}^n. \quad (2.14)$$

The corresponding factor group is commonly denoted by \mathbb{R}^n/Λ . If Λ is a lattice of maximal rank n , this factor group is called a *torus* of dimension n .

With the concept of lattices at hand, the notion of periodicity can be made precise as follows. Given a Bravais lattice $\Lambda \subset \mathbb{R}^n$ and a non-empty set Y , a function $f : \mathbb{R}^n \rightarrow Y$ is called *periodic with respect to Λ* or simply *Λ -periodic*, if it is invariant under the action of $\text{Trn}(\Lambda)$, i.e., if

$$f \circ \tau_{\mathbf{a}} = f \quad \text{for all } \mathbf{a} \in \Lambda.$$

For every Bravais lattice Λ there exist fundamental regions of the corresponding group of lattice translations $\text{Trn}(\Lambda)$. The interior of such a fundamental region is usually referred to as a *primitive domain* of Λ . The closure is referred to as a *primitive cell* of Λ . The primitive cells of a lattice $\Lambda \subset \mathbb{R}^n$ are bounded sets, if and only if Λ is a lattice of maximal rank n . All primitive cells of such a lattice Λ share the same measure. This measure is commonly referred to as the *discriminant* of Λ . A primitive cell, which is uniquely defined for every Bravais lattice, is the so-called *Wigner–Seitz cell*. Given a lattice $\Lambda \subset \mathbb{R}^3$, the Wigner–Seitz cell W_Λ of Λ is given by

$$W_\Lambda := \{ \mathbf{x} \in \mathbb{R}^n \mid |\mathbf{x}| \leq |\mathbf{x} - \mathbf{a}| \text{ for all } \mathbf{a} \in \Lambda \}. \quad (2.15)$$

The set W_Λ consists of all points, which are closer to the origin than to any other point in the lattice Λ . Hence, W_Λ is equal to the closure of the origin’s Voronoi region with respect to Λ (see e.g. Section 1.1 in [47]). Since every Voronoi region can be characterized as the intersection of finitely many open half-spaces, it follows that every Wigner–Seitz cell is a closed, convex polytope. Finally, we remark that every Wigner–Seitz cell is point-symmetric with respect to the origin. This is due to the fact that every Bravais lattice is point-symmetric with respect to each of its lattice points. For every lattice Λ we consequently have that

$$|\mathbf{x}| \leq \frac{1}{2} \text{diam}(W_\Lambda) \quad \text{for all } \mathbf{x} \in W_\Lambda.$$

Finally, we introduce the concept of reciprocal lattices. For every Bravais lattice $\Lambda \subset \mathbb{R}^n$ of rank r there exists a uniquely defined lattice $\widehat{\Lambda} \subset \mathbb{R}^n$, which is also of rank r , such that

$$\mathbf{a} \cdot \mathbf{b} \in 2\pi\mathbb{Z} \quad \text{for all } \mathbf{a} \in \Lambda, \mathbf{b} \in \widehat{\Lambda}$$

This lattice $\widehat{\Lambda}$ is called the *reciprocal lattice* of Λ . If $\mathbf{A} \in \mathbb{R}^{n \times r}$ is a matrix, whose columns form a \mathbb{Z} -linear basis of Λ , then the columns of the matrix $2\pi\mathbf{A}(\mathbf{A}^T\mathbf{A})^{-1}$ form a \mathbb{Z} -linear basis of the reciprocal lattice $\widehat{\Lambda}$. Reciprocal lattices appear e.g. in the definition of Fourier series (see Section 4.2). As will be discussed in Section 3.4, they also play an important role in the definition of so-called Bloch modes, which form a special class of solutions of differential equations with periodic coefficients.

Chapter 3

The Mathematical Model

In this chapter we develop a mathematical model, which will allow us to describe the medium structure of photonic crystals, as well as the propagation of light waves in photonic crystals. We start off by making certain assumptions on the medium structure of photonic crystals in Section 3.1. In Section 3.2 we investigate in detail the symmetries, which can arise for these medium structures. We also introduce some terminology from crystallography, since this terminology is often used in the literature on photonic crystals. In Section 3.3 we present the mathematical model, which describes the propagation of light waves in certain dielectric media. A specific ansatz for light waves in spatially periodic media is introduced in Section 3.4. This so-called Bloch ansatz leads to a family of eigenvalue problems, which determines whether or not a light wave is able to propagate inside a photonic crystal. In Section 3.5 we consider the special case of two-dimensional photonic crystals and introduce simplified models for electromagnetic wave propagation therein.

3.1 Modelling Photonic Crystals

Throughout this work we assume that a photonic crystal is a medium, which occupies the entire three-dimensional Euclidean space \mathbb{R}^3 . In doing so we focus on the propagation of light waves inside the photonic crystal, neglecting any effects occurring at its boundary. Furthermore we assume that a photonic crystal consists of a finite number of different materials. The crystal's *medium structure*, i.e. the distribution of the different materials within the crystal, can then be represented by a function $\chi : \mathbb{R}^3 \rightarrow \mathbb{N}$, in the sense that every function value $\chi(\mathbf{x})$ represents the material at the corresponding point $\mathbf{x} \in \mathbb{R}^3$. In the following we shall call this function the *crystal function*. We remark that many artificially created photonic crystals simply consist of a dielectric background material, such as silicon, into which some sort of holes are placed periodically. Such crystals can be represented

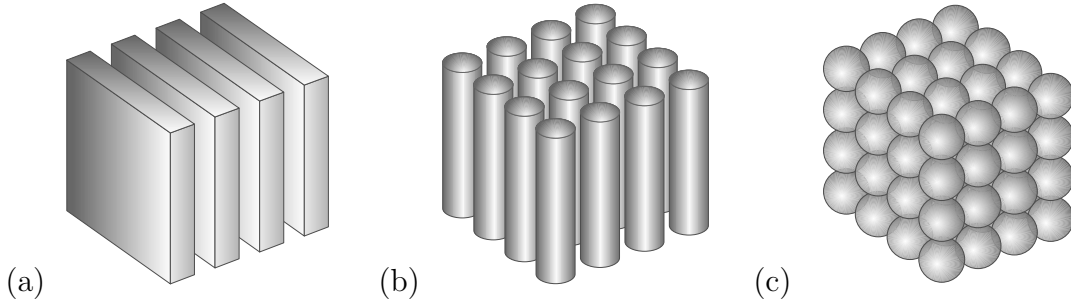


Figure 3.1: Examples of a one-dimensional (a), two-dimensional (b), and three-dimensional (c) photonic crystal.

by a two-valued crystal function with values in the set $\{1, 2\}$, where 1 models air and 2 the dielectric.

As characterizing feature, every photonic crystal exhibits some kind of spatial periodicity, which is captured by its crystal function. A photonic crystal is called *three-dimensional*, if its crystal function χ is periodic with respect to a Bravais lattice $\Lambda \subset \mathbb{R}^3$ of rank 3. Recall that such a lattice is given as the \mathbb{Z} -linear hull of three linearly independent primitive lattice vectors (see Section 2.5). The function χ is Λ -periodic, if and only if it is invariant under all lattice translations, i.e., if

$$\chi \circ \tau_{\mathbf{a}} = \chi \quad \text{for all } \mathbf{a} \in \Lambda,$$

where $\tau_{\mathbf{a}}$ denotes the translation operator defined in (2.12). A photonic crystal is called *one-dimensional* or *two-dimensional*, if its crystal function is periodic with respect to a Bravais lattice of rank 1 or 2, respectively, and if the crystal function is furthermore invariant under all translations by vectors that are perpendicular to Bravais lattice. In summary, a photonic crystal is called *r-dimensional*, $r \in \{1, 2, 3\}$, if there exists a Bravais lattice $\Lambda \subset \mathbb{R}^3$ of rank r , such that the crystal function χ satisfies

$$\chi \circ \tau_{\mathbf{y}} = \chi \quad \text{for all } \mathbf{y} \in \Lambda \oplus \Lambda^\perp.$$

Note that the \mathbb{R} -linear hull of a lattice $\Lambda \subset \mathbb{R}^3$ of rank r is a linear subspace of \mathbb{R}^3 of dimension r . Therefore, the orthogonal complement Λ^\perp of Λ is a linear subspace of \mathbb{R}^3 of dimension $3 - r$. Simply put, a photonic crystal is called three-dimensional, if the periodicity of its medium structure extends into all three space dimensions. A photonic crystal is called two-dimensional, if its medium structure is *homogeneous* in one spatial direction while being periodic in all directions, which are perpendicular to that direction. Finally, a photonic crystal is called one-dimensional if its medium structure is periodic in one spatial direction and homogeneous on every plane, which is perpendicular to that direction. Examples of one-dimensional, two-dimensional, and three-dimensional photonic crystals are depicted Figure 3.1.

Without loss of generality, we assume in the following that the Bravais lattice Λ , which describes the spatial periodicity of an r -dimensional photonic crystal, can be chosen such that

$$\text{span}(\Lambda) = \{\mathbf{x} \in \mathbb{R}^3 \mid \mathbf{x}_i = 0 \text{ for all } i > r\},$$

where $\text{span}(\Lambda)$ denotes the \mathbb{R} -linear hull of Λ . With this the lattice Bravais lattice can be identified with a lattice in \mathbb{R}^r of maximal rank r . Furthermore, the crystal function can be identified with a function mapping \mathbb{R}^r into the set of material indices. These identification allow us to treat r -dimensional photonic crystals as r -dimensional media.

3.2 Crystal Symmetries

In this section we discuss certain aspects related to the symmetry of photonic crystals. By definition, every photonic crystal exhibits some kind of periodicity, which can be described in terms of a Bravais lattice. Aside from this translational symmetry the medium structures of most photonic crystals exhibit further symmetries. Knowing these non-translational symmetries is essential when analysing the crystal's optical properties.

The scientific discipline, which is primarily concerned with the study and classification of symmetries in spatially periodic media, is *crystallography*. In crystallography, sophisticated classification systems and specialized notations are introduced in order to characterize and designate the symmetry groups of periodic media. In this section we introduce some of the crystallographic terminology, which is often used in the literature on photonic crystals. The main intent of this section is to provide a dictionary, which translates certain notions from crystallography into group theoretic concepts. A more detailed discussion can be found e.g. in [6] and [62]. Specialized notation systems, such as the Schönflies notation or the Hermann–Mauguin notation, which are commonly used in crystallography, are not introduced in this section. They are covered by most standard textbooks on crystallography, such as [15].

In Section 2.5 we introduced *primitive cells* of Bravais lattices as the closures of fundamental regions of the group of lattice translations. A related concept in crystallography is that of unit cells. Let $\Lambda \subset \mathbb{R}^n$ be a Bravais lattice of rank $r \leq n$. Then, a set $P \subset \mathbb{R}^n$ of the form

$$P := \{x_1 \mathbf{a}^{(1)} + \dots + x_r \mathbf{a}^{(r)} \mid x_1, \dots, x_r \in [0, 1]\},$$

where $\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(r)} \in \Lambda$ are linearly independent lattice vectors, is called a *unit cell* of Λ . Note that the unit cells of a lattice of rank 3 are parallelepipeds, while

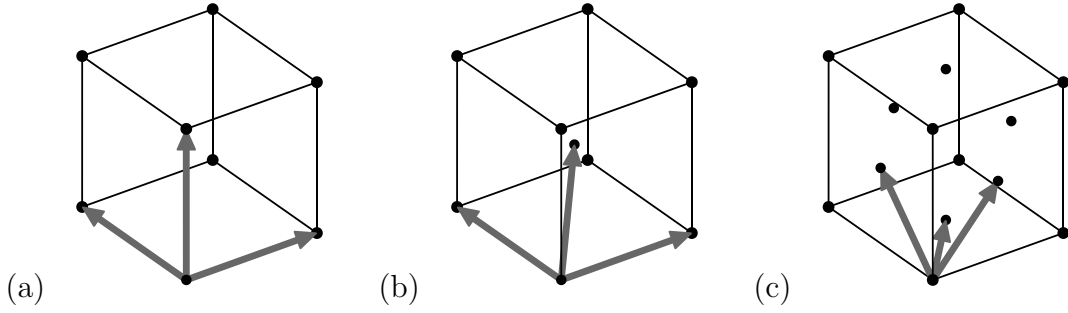


Figure 3.2: The standard unit cells of a simple cubic (a), body-centered cubic (b), and face-centered cubic (c) lattice. The arrows indicate primitive lattice vectors.

the unit cells of a lattice of rank 2 are parallelograms. A unit cell of a lattice Λ is called *primitive*, if it is spanned by primitive vectors of Λ only (see Section 2.5). Clearly, every primitive unit cell of Λ is a primitive cell of Λ . A unit cell, which is not primitive, is called *centered*.

In crystallography, Bravais lattices are categorized according to the shape of their primitive unit cells. For every lattice type, crystallography defines a so-called *standard unit cell*. For many lattice types these standard unit cells are primitive or centered unit cells in the sense that was defined above. For some lattice types, however, other polytopes are considered as standard unit cells. In three space dimensions there exist 14 different *crystallographic lattice types* in total. Each lattice type belongs to exactly one of 7 so-called *crystallographic lattice systems*. In this work, we only consider the three lattice types, which constitute to the so-called *cubic* lattice system. These ones are the so-called simple cubic, body-centered cubic (bcc), and face-centered cubic (fcc) lattices.

A Bravais lattice Λ of rank 3 is called *simple cubic*, if there exist primitive vectors $\mathbf{a}^{(1)}, \mathbf{a}^{(2)}, \mathbf{a}^{(3)}$ of Λ that satisfy

$$\begin{aligned} |\mathbf{a}^{(1)}| &= |\mathbf{a}^{(2)}| = |\mathbf{a}^{(3)}| =: a > 0, \\ \angle(\mathbf{a}^{(1)}, \mathbf{a}^{(2)}) &= \angle(\mathbf{a}^{(2)}, \mathbf{a}^{(3)}) = \angle(\mathbf{a}^{(1)}, \mathbf{a}^{(3)}) = \frac{\pi}{2}. \end{aligned}$$

Clearly, a Bravais lattice is a simple cubic lattice, if and only if its primitive unit cells are cubes. The length a of the cubes' edges is called the *lattice constant*. Note that every simple cubic lattice is similar to the *standard simple cubic lattice* \mathbb{Z}^3 . Choosing $\mathbf{e}^{(1)}, \mathbf{e}^{(2)},$ and $\mathbf{e}^{(3)}$ as primitive vectors, where $\mathbf{e}^{(i)}$ denotes the i -th standard basis vector in \mathbb{R}^3 for $i = 1, 2, 3$, one obtains the primitive unit cell P_3 of \mathbb{Z}^3 , which is given by

$$P_3 := [0, 1]^3. \quad (3.1)$$

In crystallography, P_3 is the standard unit cell of \mathbb{Z}^3 .

Next, we introduce body-centered cubic lattices. A Bravais lattice Λ of rank 3 is called *body-centered cubic* or *bcc*, if there exist primitive vectors $\mathbf{a}^{(1)}, \mathbf{a}^{(2)}, \mathbf{a}^{(3)}$ of Λ , such that

$$\begin{aligned} |\mathbf{a}^{(1)}| &= |\mathbf{a}^{(2)}| =: a > 0, & |\mathbf{a}^{(3)}| &= \frac{\sqrt{3}a}{2}, \\ \angle(\mathbf{a}^{(1)}, \mathbf{a}^{(2)}) &= \frac{\pi}{2}, & \angle(\mathbf{a}^{(2)}, \mathbf{a}^{(3)}) &= \angle(\mathbf{a}^{(1)}, \mathbf{a}^{(3)}) = \frac{\pi}{4}. \end{aligned}$$

The primitive cells of a bcc lattice are parallelepipeds with exactly two parallel square faces. The edge lengths of these faces are given by the lattice constant a . The distance between the square faces is equal to $a/2$. One quickly discovers that every bcc lattice can be characterized as the union of two staggered simple cubic lattices. More precisely, every bcc lattice is similar to the *standard bcc lattice* Λ_{bcc} , which is defined as

$$\Lambda_{\text{bcc}} := \left\{ z_1 \mathbf{e}^{(1)} + z_2 \mathbf{e}^{(2)} + \frac{z_3}{2} (\mathbf{e}^{(1)} + \mathbf{e}^{(2)} + \mathbf{e}^{(3)}) \mid z_1, z_2, z_3 \in \mathbb{Z} \right\}. \quad (3.2)$$

One easily verifies that \mathbb{Z}^3 is a sublattice of Λ_{bcc} (see Section 2.5). Because of this, the standard unit cell P_3 of \mathbb{Z}^3 is a centered unit cell of Λ_{bcc} . In crystallography, P_3 is considered to be the standard unit cell of Λ_{bcc} .

The third class of Bravais lattices, which belong to the cubic lattice system, consists of the so-called face-centered cubic lattices. A Bravais lattice Λ of rank 3 is called *face-centered cubic* or *fcc*, if there exist primitive vectors $\mathbf{a}^{(1)}, \mathbf{a}^{(2)}, \mathbf{a}^{(3)}$ of Λ , such that

$$\begin{aligned} |\mathbf{a}^{(1)}| &= |\mathbf{a}^{(2)}| = |\mathbf{a}^{(3)}| := \frac{\sqrt{2}a}{2}, & a &> 0, \\ \angle(\mathbf{a}^{(1)}, \mathbf{a}^{(2)}) &= \angle(\mathbf{a}^{(2)}, \mathbf{a}^{(3)}) = \angle(\mathbf{a}^{(1)}, \mathbf{a}^{(3)}) = \frac{\pi}{3}. \end{aligned}$$

The primitive unit cells of a fcc lattice are rhombohedra. The edge lengths of these rhombohedra is given by $\sqrt{2}a/2$, where a denotes the lattice constant. Every fcc lattice is therefore similar to the *standard fcc lattice* Λ_{fcc} , which is defined by

$$\Lambda_{\text{fcc}} := \left\{ \frac{z_1}{2} (\mathbf{e}^{(1)} + \mathbf{e}^{(2)}) + \frac{z_2}{2} (\mathbf{e}^{(1)} + \mathbf{e}^{(3)}) + \frac{z_3}{2} (\mathbf{e}^{(2)} + \mathbf{e}^{(3)}) \mid z_1, z_2, z_3 \in \mathbb{Z} \right\}. \quad (3.3)$$

Since \mathbb{Z}^3 is also a sublattice of Λ_{fcc} , the cube P_3 is also considered to be the standard unit cell of Λ_{fcc} . In Figure 3.2 we depict the standard unit cells of all three lattice types, which belong to the cubic lattice system.

In two space dimensions, crystallography distinguishes 5 different lattice types. Here, we only consider square and hexagonal lattices. A Bravais lattice Λ of rank 2 is called *square*, if there exist primitive vectors $\mathbf{a}^{(1)}, \mathbf{a}^{(2)}$ of Λ , such that

$$|\mathbf{a}^{(1)}| = |\mathbf{a}^{(2)}| =: a > 0,$$

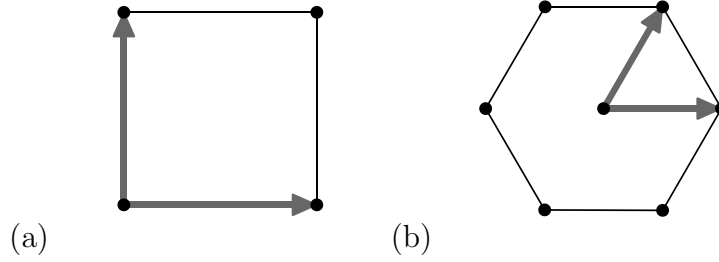


Figure 3.3: The standard unit cells of a square (a) and a hexagonal lattice. The arrows indicate primitive lattice vectors.

$$\angle(\mathbf{a}^{(1)}, \mathbf{a}^{(2)}) = \frac{\pi}{2}.$$

As the name indicates, square lattices are exactly those Bravais lattices, which feature square primitive unit cells. The side lengths of these squares are given by the lattice constant a . Every square lattice is similar to the *standard square lattice* \mathbb{Z}^2 . The standard unit cell of \mathbb{Z}^2 is the square P_2 , which is defined by

$$P_2 := [0, 1]^2. \quad (3.4)$$

Clearly, P_2 is a primitive unit cell of \mathbb{Z}^2 .

Finally, we introduce the so-called hexagonal lattices of rank 2. A Bravais lattice Λ of rank 2 is called *hexagonal*, if one can find primitive vectors $\mathbf{a}^{(1)}, \mathbf{a}^{(2)}$ of Λ , such that

$$\begin{aligned} |\mathbf{a}^{(1)}| &= |\mathbf{a}^{(2)}| =: a > 0, \\ \angle(\mathbf{a}^{(1)}, \mathbf{a}^{(2)}) &= \frac{\pi}{3}. \end{aligned}$$

The primitive unit cells of a hexagonal lattice are rhombi with side lengths equal to the lattice constant a and angles measuring 60° and 120° . Every hexagonal lattice is hence similar to the *standard hexagonal lattice* Λ_{hex} , given by

$$\Lambda_{\text{hex}} := \left\{ z_1 \mathbf{e}^{(1)} + z_2 \left(\frac{1}{2} \mathbf{e}^{(1)} + \frac{\sqrt{3}}{2} \mathbf{e}^{(2)} \right) \mid z_1, z_2 \in \mathbb{Z} \right\}, \quad (3.5)$$

where $\mathbf{e}^{(1)}$ and $\mathbf{e}^{(2)}$ denote the standard basis vectors of \mathbb{R}^2 . The name *hexagonal* alludes to the fact that all lattice points in a hexagonal lattice are equidistant to their nearest neighbours. Hence, by taking the convex hull of the nearest neighbours of a lattice point, one obtains a regular hexagon. The hexagon with the origin at its center is defined as the standard unit cell of Λ_{hex} . In Figure 3.3 we

depict the standard unit cells of square and hexagonal lattices. It should be noted that there are also Bravais lattices of rank 3, which are called hexagonal. Such lattices are obtained through periodic translations of a hexagonal lattice of rank 2 in along the direction, which is perpendicular to that lattice.

Recall that the medium structure of an r -dimensional photonic crystal can be described by a crystal function $\chi : \mathbb{R}^r \rightarrow \mathbb{N}$ (see Section 3.1). The periodicity of the medium structure is represented by a Bravais lattice Λ of rank r , with respect to which χ is Λ -periodic. In the following, we assume that the lattice Λ is chosen maximally, i.e., we assume that there are no translations under which χ is invariant apart from those contained in $\text{Trn}(\Lambda)$. Recall that $\text{Trn}(\Lambda)$ denotes the group of translations, which are generated by Λ (see Section 2.5). One can easily verify that $\text{Trn}(\Lambda)$ is a normal subgroup of $\text{Iso}(\mathbb{R}^r)_\chi$, the symmetry group of χ . In contrast to $\text{Trn}(\Lambda)$, which only represents the translational symmetry of the photonic crystal's medium structure, $\text{Iso}(\mathbb{R}^r)_\chi$ represents its complete symmetry. Therefore, $\text{Iso}(\mathbb{R}^r)_\chi$ can contain symmetry operations, which correspond to non-translational symmetries arising from the way the different materials of the photonic crystal are arranged within a primitive cell of Λ . Note, however, that $\text{Iso}(\mathbb{R}^r)_\chi$ is always a subgroup of $\text{Iso}(\mathbb{R}^r)_\Lambda$, the symmetry group of Λ , since every symmetry operation of χ in particular must take Λ into itself. In summary, we therefore have

$$\text{Trn}(\Lambda) \trianglelefteq \text{Iso}(\mathbb{R}^r)_\chi \leq \text{Iso}(\mathbb{R}^r)_\Lambda,$$

where \trianglelefteq indicates the normal subgroup relation.

In crystallography, the complete symmetry of a periodic medium is characterized in terms of so-called *crystallographic space groups*. The complete symmetry of an r -dimensional photonic crystal's medium structure is said to be given by the space group G_{space} , if there exists a crystal function χ , such that its symmetry group $\text{Iso}(\mathbb{R}^r)_\chi$ is isomorphic to G_{space} . It turns out that there exist exactly 230 different crystallographic space groups in three space dimensions. This means that the complete symmetry of the medium structure of every three-dimensional photonic crystal is given by exactly one of these 230 groups. It is further known that in two space dimensions there exists a total of 17 different space groups, which are sometimes referred to as the 17 *wallpaper groups*.

As was stated before, space groups represent the complete symmetry of a periodic medium. The complete symmetry of such a medium usually consists of some sort of translational symmetry, which is represented by a Bravais lattice, as well as of other, non-translational symmetries. In crystallography, the non-translational symmetries of a periodic medium are characterized in terms of so-called *crystallographic point groups*. The non-translational symmetries of an r -dimensional photonic crystal's medium structure, whose complete symmetry is given by the space group G_{space} and whose translational symmetry is given by a Bravais lattice

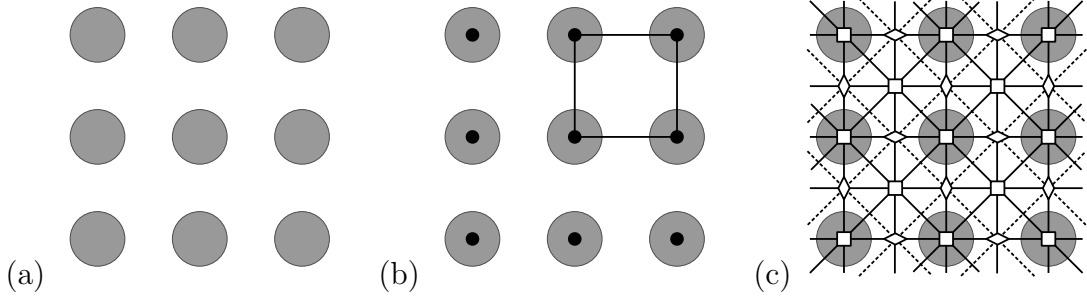


Figure 3.4: Example of a two-dimensional photonic crystal's medium structure (a) with its translational (b) and non-translational symmetries (c). The space group is isomorphic to $\mathbb{Z}^2 \rtimes O_2(\mathbb{Z})$.

Λ , are said to be given by the point group G_{point} , if there exists a crystal function χ , such that $\text{Iso}(\mathbb{R}^r)_\chi$ is isomorphic to G_{space} and such that $\text{Iso}(\mathbb{R}^r)_\chi / \text{Trn}(\Lambda)$ is isomorphic to G_{point} . It turns out that there exist exactly 32 different crystallographic point groups in three space dimensions. In two space dimensions there exist exactly 10 different point groups.

In the following, we present two examples, which illustrate the concept of point groups and space groups in two space dimensions. First, we consider a two-dimensional photonic crystal, whose medium structure is depicted in Figure 3.4(a). Without loss of generality, we can assume that the medium structure is represented by a crystal function χ , which is periodic with respect to the standard square lattice \mathbb{Z}^2 . Figure 3.4(b) depicts the structure's translational symmetry. Lattice points are indicated by black dots (•), the standard unit cell is indicated by a square. Apart from the translational symmetry the medium structure also exhibits a number of *rotational symmetries*, *reflection symmetries* and *glide-reflection symmetries*. In Figure 3.4(c) these symmetries are indicated by graphical symbols. Squares (□) indicate centers of rotations by $\pi/2$, π , and $3\pi/2$ which take the medium structure into itself. Such rotation centers are called *four-fold centers of rotation*. Centers of rotations by π that take the medium structure into itself are called *two-fold centers of rotation* and indicated by lozenges (◇). Solid lines indicate *axes of reflections*, and dashed lines indicate *axis of glide-reflections* that take the medium structure into itself. A glide-reflection is defined as a reflection, which is followed by a translation along the axis of reflection. One finds that the symmetry group of χ is given by

$$\text{Iso}(\mathbb{R}^2)_\chi = \{ \varphi : \mathbb{R}^2 \rightarrow \mathbb{R}^2 \mid \varphi(\mathbf{x}) = \Theta \mathbf{x} + \mathbf{a}, \Theta \in O_2(\mathbb{Z}), \mathbf{a} \in \mathbb{Z}^2 \},$$

where $O_2(\mathbb{Z})$ denotes the group of orthogonal 2×2 -matrices with integer components. Furthermore, one can easily show that $\text{Iso}(\mathbb{R}^2)_\chi$ is isomorphic to the

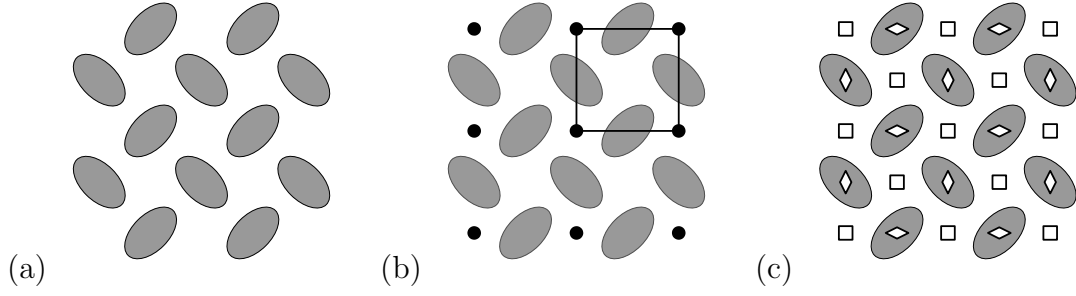


Figure 3.5: Example of a two-dimensional photonic crystal's medium structure (a) with its translational (b) and non-translational symmetries (c). The space group is isomorphic to $\mathbb{Z}^2 \rtimes \text{SO}_2(\mathbb{Z})$.

standard semidirect product $\mathbb{Z}^2 \rtimes \text{O}_2(\mathbb{Z})$. Recall that $\mathbb{Z}^2 \rtimes \text{O}_2(\mathbb{Z})$ is defined as the Cartesian product of \mathbb{Z}^2 and $\text{O}_2(\mathbb{Z})$ endowed with the group operation, which is given by

$$(\mathbf{a}^{(1)}, \Theta^{(1)})(\mathbf{a}^{(2)}, \Theta^{(2)}) := (\mathbf{a}^{(1)} + \Theta^{(1)}\mathbf{a}^{(2)}, \Theta^{(1)}\Theta^{(2)})$$

for all $\mathbf{a}^{(1)}, \mathbf{a}^{(2)} \in \mathbb{Z}^2$ and all $\Theta^{(1)}, \Theta^{(2)} \in \text{O}_2(\mathbb{Z})$. The space group, which corresponds to $\text{Iso}(\mathbb{R}^2)_\chi$, is the wallpaper group p4mm (cf. Chapter 26 in [6]). The space group of the medium structure is isomorphic to $\text{O}_2(\mathbb{Z})$. We remark that the three-dimensional analogue of the medium structure depicted in Figure 3.4(a) is a \mathbb{Z}^3 -periodic arrangement of balls. The corresponding space group, which is denoted by $\text{Pm}\bar{3}\text{m}$ (cf. Chapter 10 in [15]), is isomorphic to $\mathbb{Z}^3 \rtimes \text{O}_3(\mathbb{Z})$, the point group being isomorphic to $\text{O}_3(\mathbb{Z})$.

As a second example, we briefly consider a two-dimensional photonic crystal, whose medium structure is depicted by Figure 3.5(a). As in the previous example, the translational symmetry is given by the standard square lattice \mathbb{Z}^2 . The lattice and its standard unit cell are depicted in Figure 3.5(b). The medium structure exhibits the same rotational symmetries as the structure depicted in Figure 3.4(a). However, it lacks any reflection symmetries. Thus, we do not find any solid or dashed lines in Figure 3.5(c). The symmetry group of the crystal function χ is given by

$$\text{Iso}(\mathbb{R}^2)_\chi = \{\varphi : \mathbb{R}^2 \rightarrow \mathbb{R}^2 \mid \varphi(\mathbf{x}) = \Theta\mathbf{x} + \mathbf{a}, \Theta \in \text{SO}_2(\mathbb{Z}), \mathbf{a} \in \mathbb{Z}^2\}.$$

Here as in the following, $\text{SO}_2(\mathbb{Z})$ denotes the group of orthogonal 2×2 -matrices, whose determinant equals 1. The symmetry group is isomorphic to $\mathbb{Z}^2 \rtimes \text{SO}_2(\mathbb{Z})$, and the corresponding space group is denoted by p4 (cf. Chapter 26 in [6]). The point group of the medium structure is isomorphic to $\text{SO}_2(\mathbb{Z})$.

We conclude with the remark that the representation of point groups by subgroups of $\text{O}_r(\mathbb{R})$, r being the dimension of the photonic crystal, is particularly

useful, when symmetrization operations have to be realized. Such operations will play an important role in the Chapters 7 and 8.

3.3 Wave Propagation in Linear Dielectrics

In this section we introduce a mathematical model for the propagation of light waves in certain dielectric media. This model will serve as a basis for more specific models, which are developed in Sections 3.4 and 3.5.

The standard theory, which describes the propagation of electromagnetic waves, such as light waves, is *classical electromagnetism*. Central to this theory are Maxwell's equations, a system consisting of two pairs of coupled partial differential equations, which capture fundamental properties of electromagnetic waves. The system of Maxwell's equations alone, however, does not determine the propagation of an electromagnetic wave completely. Further constitutive relations, which describe the wave's interaction with the ambient medium, have to be introduced in order to close the system. For each set of constitutive relations one obtains a different mathematical model for electromagnetic wave propagation. Here, we present the standard model for time-harmonic wave propagation in non-magnetic, non-dispersive, non-absorptive, linear, isotropic dielectrics, which is commonly used for photonic crystals (see e.g. [42], [43], [60], [66]). A detailed derivation of this model can be found in most standard textbooks on classical electromagnetism, such as [41], or in the original work of Maxwell [55]. Here, as in the remaining sections of this chapter, we focus on a formal derivation of the relevant formulas, leaving aside questions concerning the existence and uniqueness of solutions.

In classical electromagnetism electromagnetic waves are described by four vector fields $\mathbf{E}, \mathbf{H}, \mathbf{D}, \mathbf{B} : \mathbb{R} \times \mathbb{R}^3 \rightarrow \mathbb{R}^3$, which are solutions of *Maxwell's equations*. Each vector field is a function in the variables (t, \mathbf{x}) , where $t \in \mathbb{R}$ represents a point in time and $\mathbf{x} \in \mathbb{R}^3$ a position in space. The vector fields \mathbf{E} and \mathbf{H} represent the wave's *electric* and *magnetic field*, respectively. The vector field \mathbf{D} represents the wave's *electric displacement field*, and the vector field \mathbf{B} the wave's *magnetic flux density*. The fields \mathbf{D} and \mathbf{B} depend on the medium, in which the electromagnetic wave propagates, while the fields \mathbf{E} and \mathbf{H} describe the wave in free space. Under the assumption that all physical quantities are given in SI units (cf. [18]), the so-called macroscopic form of Maxwell's equations reads

$$\left\{ \begin{array}{ll} \nabla \times \mathbf{E} + \dot{\mathbf{B}} = \mathbf{0} & \text{in } \mathbb{R} \times \mathbb{R}^3 \quad (\text{Faraday's law}) \\ \nabla \times \mathbf{H} - \dot{\mathbf{D}} = \mathbf{J}_f & \text{in } \mathbb{R} \times \mathbb{R}^3 \quad (\text{Ampère's law}) \\ \nabla \cdot \mathbf{D} = \rho_f & \text{in } \mathbb{R} \times \mathbb{R}^3 \quad (\text{Gauss's law}) \\ \nabla \cdot \mathbf{B} = 0 & \text{in } \mathbb{R} \times \mathbb{R}^3 \quad (\text{Gauss's law for magnetism}) \end{array} \right. \quad (3.6)$$

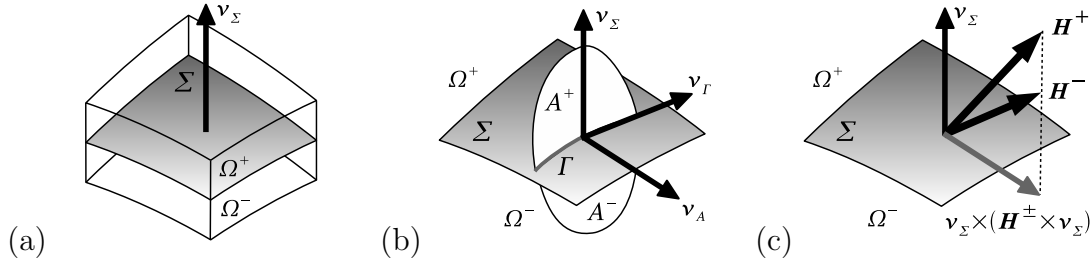


Figure 3.6: The layers Ω^+ and Ω^- are separated by the surface Σ (a). The surface $A = A^+ \cup A^-$ intersects Σ perpendicularly in Γ (b). In general, only the tangential component of the magnetic field \mathbf{H} is continuous at Σ (c).

Here, $\dot{\mathbf{D}}$ and $\dot{\mathbf{B}}$ denote the partial derivatives of the vector fields \mathbf{D} and \mathbf{B} with respect to the time variable t . The curl and divergence operators in (3.6) are understood to act on the vector fields with respect to the space variable \mathbf{x} . The vector field $\mathbf{J}_f : \mathbb{R} \times \mathbb{R}^3 \rightarrow \mathbb{R}^3$, as well as the scalar field $\rho_f : \mathbb{R} \times \mathbb{R}^3 \rightarrow \mathbb{R}$ are given functions, which model the *free current density* and the *free charge density* inside the medium. In this work we only consider media, in which neither free currents nor free charges exist. As a consequence, we have

$$\mathbf{J}_f = \mathbf{0}, \quad (3.7)$$

$$\rho_f = 0. \quad (3.8)$$

From top to bottom, the equations in (3.6) are known as *Faraday's law*, *Ampère's law*, *Gauss' law*, and *Gauss' law for magnetism*. The latter one, in particular, states that there are no magnetic monopoles.

Ampère's law implies an important property of the magnetic field \mathbf{H} concerning its transition at an interface. To show this, let us assume that Σ is smooth surface in \mathbb{R}^3 , which is oriented by a *unit normal field* $\boldsymbol{\nu}_\Sigma : \Sigma \rightarrow \mathbb{S}^2$, where \mathbb{S}^2 denotes the three-dimensional unit sphere. Such a surface could be, for example, the interface between two different media. Given some positive number $\eta > 0$, we define the layers $\Omega^+, \Omega^- \subset \mathbb{R}^3$ by

$$\Omega^+ := \{ \mathbf{x} + \delta \boldsymbol{\nu}_\Sigma(\mathbf{x}) \mid \mathbf{x} \in \Sigma, 0 \leq \delta \leq \eta \},$$

$$\Omega^- := \{ \mathbf{x} - \delta \boldsymbol{\nu}_\Sigma(\mathbf{x}) \mid \mathbf{x} \in \Sigma, 0 \leq \delta \leq \eta \},$$

as well as the set $\Omega := \Omega^+ \cup \Omega^-$. Note that the intersection of Ω^+ and Ω^- coincides with Σ . For illustration, see Figure 3.6(a). Given a vector field $\mathbf{F} : \mathbb{R} \times \mathbb{R}^3 \rightarrow \mathbb{R}^3$, we define $\mathbf{F}^+ : \mathbb{R} \times \Omega^+ \rightarrow \mathbb{R}^3$ and $\mathbf{F}^- : \mathbb{R} \times \Omega^- \rightarrow \mathbb{R}^3$ by

$$\mathbf{F}^+(t, \mathbf{x}) := \begin{cases} \lim_{\delta \rightarrow 0^+} \mathbf{F}(t, \mathbf{x} + \delta \boldsymbol{\nu}_\Sigma(\mathbf{x})) & \text{if } \mathbf{x} \in \Sigma, \\ \mathbf{F}(t, \mathbf{x}) & \text{else} \end{cases} \quad \text{for all } t \in \mathbb{R}, \mathbf{x} \in \Omega^+, \quad (3.9)$$

$$\mathbf{F}^-(t, \mathbf{x}) := \begin{cases} \lim_{\delta \rightarrow 0^+} \mathbf{F}(t, \mathbf{x} - \delta \boldsymbol{\nu}_\Sigma(\mathbf{x})) & \text{if } \mathbf{x} \in \Sigma, \\ \mathbf{F}(t, \mathbf{x}) & \text{else} \end{cases} \quad \text{for all } t \in \mathbb{R}, \mathbf{x} \in \Omega^-. \quad (3.10)$$

Let us now assume that $\Gamma \subset \Sigma$ is a smooth curve, which is oriented by a *unit velocity field* $\boldsymbol{\nu}_\Gamma : \Gamma \rightarrow \mathbb{S}^2$. Assuming that Γ is neither closed nor self-intersecting, one can find a simply connected, bounded, smooth surface $A \subset \mathbb{R}^3$, which is contained in the set Ω , and which intersects the surface Σ exactly in Γ . We further assume that the sets $A^+ := \Omega^+ \cap A$ and $A^- := \Omega^- \cap A$ are non-empty, and that the surface A is oriented by a unit normal field $\boldsymbol{\nu}_A : A \rightarrow \mathbb{S}^2$ (see Figure 3.6(b)). Now, let $t \in \mathbb{R}$ be an arbitrary point in time. Then, we obtain by Stokes' theorem and Ampère's law that

$$\begin{aligned} & \int_\Gamma \boldsymbol{\nu}_\Gamma \cdot (\mathbf{H}^+(t, \cdot) - \mathbf{H}^-(t, \cdot)) \\ &= \int_{\partial A^+} \boldsymbol{\nu}_{\partial A^+} \cdot \mathbf{H}^+(t, \cdot) + \int_{\partial A^-} \boldsymbol{\nu}_{\partial A^-} \cdot \mathbf{H}^-(t, \cdot) - \int_{\partial A} \boldsymbol{\nu}_{\partial A} \cdot \mathbf{H}(t, \cdot) \\ &= \int_{A^+} \boldsymbol{\nu}_A \cdot \nabla \times \mathbf{H}^+(t, \cdot) + \int_{A^-} \boldsymbol{\nu}_A \cdot \nabla \times \mathbf{H}^-(t, \cdot) - \int_A \boldsymbol{\nu}_A \cdot \nabla \times \mathbf{H}(t, \cdot) \\ &= \int_{A^+} \boldsymbol{\nu}_A \cdot \dot{\mathbf{D}}^+(t, \cdot) + \int_{A^-} \boldsymbol{\nu}_A \cdot \dot{\mathbf{D}}^-(t, \cdot) - \int_A \boldsymbol{\nu}_A \cdot \dot{\mathbf{D}}(t, \cdot) \\ &= 0. \end{aligned} \quad (3.11)$$

Here $\boldsymbol{\nu}_{\partial A^+}$, $\boldsymbol{\nu}_{\partial A^-}$, and $\boldsymbol{\nu}_{\partial A}$ denote the unit velocity fields on the boundary curves ∂A^+ , ∂A^- , and ∂A of A^+ , A^- , and A , respectively. Note that the identity (3.11) only holds under the assumption stated in (3.7), namely that there are no free currents, and under the assumption that the vector field $\dot{\mathbf{B}}$ is smooth up to the surface Σ . Since (3.11) can be established for every oriented, smooth curve Γ on Σ , we conclude that

$$\boldsymbol{\tau} \cdot (\mathbf{H}^+ - \mathbf{H}^-)|_\Sigma = 0 \quad (3.12)$$

holds in the sense of integrals along smooth curves on Σ for every unit vector field $\boldsymbol{\tau} : \Sigma \rightarrow \mathbb{S}^2$ that is *tangential* to Σ , i.e., that satisfies

$$\boldsymbol{\tau} \cdot \boldsymbol{\nu}_\Sigma = 0. \quad (3.13)$$

It follows from (3.12) that the vector field $(\mathbf{H}^+ - \mathbf{H}^-)|_\Sigma$ is *normal* on Σ , which means that

$$(\mathbf{H}^+ - \mathbf{H}^-)|_\Sigma = (\boldsymbol{\nu}_\Sigma \cdot (\mathbf{H}^+ - \mathbf{H}^-)|_\Sigma) \boldsymbol{\nu}_\Sigma.$$

Hence, one can easily verify that

$$\boldsymbol{\nu}_\Sigma \times (\mathbf{H}^+ - \mathbf{H}^-)|_\Sigma = \mathbf{0},$$

and that

$$\boldsymbol{\nu}_\Sigma \times ((\mathbf{H}^+ - \mathbf{H}^-)|_\Sigma \times \boldsymbol{\nu}_\Sigma) = \mathbf{0}.$$

The latter identity implies, that the transition of the so-called *tangential component* $\boldsymbol{\nu} \times (\mathbf{H} \times \boldsymbol{\nu})$ of the magnetic field \mathbf{H} is continuous at the interface Σ . Here, $\boldsymbol{\nu}$ denotes a continuous extension of the unit normal field $\boldsymbol{\nu}_\Sigma$ to Ω . As we will see below, the so-called *normal component* $\boldsymbol{\nu} \cdot \mathbf{H}$ of \mathbf{H} can be discontinuous at Σ . Figure 3.6(c) illustrates the generic transitional behavior of the magnetic field at an interface.

We now return to the system of Maxwell's equations (3.6). As one can see, this system is not closed. In order to close the system, one needs to introduce constitutive relations, which relate the fields \mathbf{D} and \mathbf{B} to the fields \mathbf{E} and \mathbf{H} . Such constitutive relations model the response of the ambient medium to the electromagnetic wave. In this work, we only consider non-dispersive, non-absorptive, linear, and isotropic dielectrics. For these media the constitutive relations are of the form

$$\mathbf{D} = \varepsilon_0 \varepsilon_r \mathbf{E}, \quad (3.14)$$

$$\mathbf{B} = \mu_0 \mu_r \mathbf{H}, \quad (3.15)$$

where $\varepsilon_r, \mu_r : \mathbb{R}^3 \rightarrow \mathbb{R}_{>0}$ are scalar fields that only depend on the space variable \mathbf{x} . Typically, ε_r and μ_r are positive, bounded, and uniformly bounded away from zero. The real numbers $\varepsilon_0, \mu_0 \in \mathbb{R}_{>0}$ are physical constants, which are called the *vacuum electric permittivity* and the *vacuum magnetic permeability*. For these constants the relation

$$\varepsilon_0 \mu_0 = \frac{1}{c_0^2} \quad (3.16)$$

holds, where c_0 denotes the *vacuum speed of light*. The scalar fields ε_r and μ_r represent the medium's *relative electric permittivity* and *relative magnetic permeability*. Both fields ε_r and μ_r determine the medium's optical density, which is measured by the space-dependent *refractive index* $n : \mathbb{R}^3 \rightarrow \mathbb{R}$. The refractive index is given by

$$n := \sqrt{\varepsilon_r \mu_r}. \quad (3.17)$$

We remark that for *non-isotropic* media the relative electric permittivity, as well as the magnetic permeability are represented by matrix-valued functions. For *dispersive*, linear media the constitutive relations between \mathbf{D} and \mathbf{E} , as well as between \mathbf{B} and \mathbf{H} are given in terms of convolutions (see e.g. Section 7.10 in [41]). Here, "dispersive" means that the properties of the medium depend on the phase frequency of the wave.

Let us, once again, consider the transitional behaviour of the magnetic field \mathbf{H} at an interface Σ . Under the above assumptions, let $\Delta \subset \Sigma$ be an oriented surface

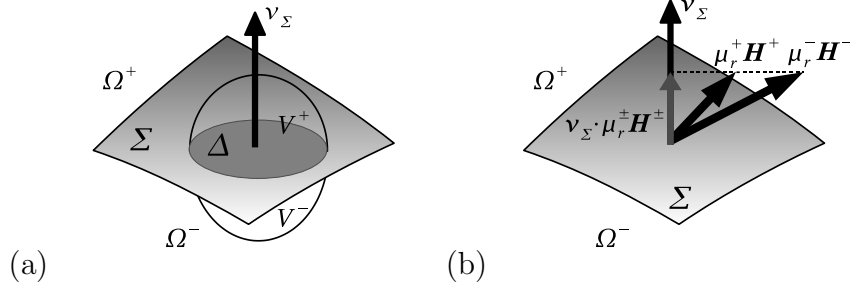


Figure 3.7: The volume $V = V^+ \cup V^-$ intersects the surface Σ in Δ (a). In general, only the normal component of $\mu_r \mathbf{H}$ is continuous at Σ (b).

on Σ , whose unit normal field is identical to that of Σ . Furthermore, let $V \subset \mathbb{R}^3$ be a smoothly bounded volume, which is contained in Ω and which intersects Σ exactly in Δ , such that $V^+ := V \cap \Omega^+$ and $V^- := V \cap \Omega^-$ are nonempty sets. For illustration, see Figure 3.7. Letting $\mu_r^+ : \Omega^+ \rightarrow \mathbb{R}$ and $\mu_r^- : \Omega^- \rightarrow \mathbb{R}$ be defined in analogy to (3.9) and (3.10), the constitutive relation (3.15), Gauss' theorem, and Gauss' law for magnetism imply that

$$\begin{aligned}
 & \int_{\Delta} \boldsymbol{\nu}_{\Sigma} \cdot (\mu_r^+ \mathbf{H}^+(t, \cdot) - \mu_r^- \mathbf{H}^-(t, \cdot)) \\
 &= \int_{\partial V^+} \boldsymbol{\nu}_{\partial V^+} \cdot \mathbf{B}^+(t, \cdot) + \int_{\partial V^-} \boldsymbol{\nu}_{\partial V^-} \cdot \mathbf{B}^-(t, \cdot) - \int_{\partial V} \boldsymbol{\nu}_{\partial V} \cdot \mu_r \mathbf{H}(t, \cdot) \\
 &= \int_{V^+} \nabla \cdot [\mathbf{B}^+(t, \cdot)] + \int_{V^-} \nabla \cdot [\mathbf{B}^-(t, \cdot)] - \int_V \nabla \cdot [\mathbf{B}(t, \cdot)] \\
 &= 0.
 \end{aligned} \tag{3.18}$$

Here, $\boldsymbol{\nu}_{\partial V^+}$, $\boldsymbol{\nu}_{\partial V^-}$, and $\boldsymbol{\nu}_{\partial V}$ denote the outer unit normal field on the boundaries ∂V^+ , ∂V^- , and ∂V of the volumes V^+ , V^- , and V , respectively. Since (3.18) can be established for every surface Δ , we conclude that

$$\boldsymbol{\nu}_{\Sigma} \cdot (\mu_r^+ \mathbf{H}^+ - \mu_r^- \mathbf{H}^-)|_{\Sigma} = 0 \tag{3.19}$$

in the sense of integrals over surfaces on Σ . It hence follows that the so-called *normal component* $\boldsymbol{\nu} \cdot \mu_r \mathbf{H}$ of the vector field $\mu_r \mathbf{H}$ is continuous across interfaces. As before, we denote by $\boldsymbol{\nu}$ a continuous extension of the unit normal field $\boldsymbol{\nu}_{\Sigma}$ to Ω . This implies, of course, that the normal component $\boldsymbol{\nu} \cdot \mathbf{H}$ of the magnetic field \mathbf{H} must be discontinuous at interfaces, where the relative magnetic permeability μ_r changes discontinuously. Figure 3.7(b) depicts the generic behaviour of the vector field $\mu_r \mathbf{H}$.

Under the assumptions (3.7), (3.8), (3.14), and (3.15), Maxwell's equations become a fully coupled system of first-order partial differential equations, which reads

$$\begin{cases} \nabla \times \mathbf{E} + \mu_0 \mu_r \dot{\mathbf{H}} = 0 & \text{in } \mathbb{R} \times \mathbb{R}^3, \\ \nabla \times \mathbf{H} - \varepsilon_0 \varepsilon_r \dot{\mathbf{E}} = 0 & \text{in } \mathbb{R} \times \mathbb{R}^3, \\ \nabla \cdot [\varepsilon_r \mathbf{E}] = 0 & \text{in } \mathbb{R} \times \mathbb{R}^3, \\ \nabla \cdot [\mu_r \mathbf{H}] = 0 & \text{in } \mathbb{R} \times \mathbb{R}^3. \end{cases} \quad (3.20)$$

One can extract a second-order system for the magnetic field \mathbf{H} from (3.20). This can be done by differentiating Faraday's law with respect to the time variable t , and by then substituting the identity

$$\dot{\mathbf{E}} = \frac{1}{\varepsilon_0 \varepsilon_r} \nabla \times \mathbf{H}, \quad (3.21)$$

which follows from Ampère's law, into the resulting equation. Together with (3.16) and Gauss' law for magnetism this yields

$$\begin{cases} \nabla \times \left[\frac{1}{\varepsilon_r} \nabla \times \mathbf{H} \right] + \frac{1}{c_0^2} \mu_r \ddot{\mathbf{H}} = 0 & \text{in } \mathbb{R} \times \mathbb{R}^3, \\ \nabla \cdot [\mu_r \mathbf{H}] = 0 & \text{in } \mathbb{R} \times \mathbb{R}^3. \end{cases} \quad (3.22)$$

Intent on studying the propagation of electromagnetic waves with a single *phase frequency* $\omega \in \mathbb{R}$, we make a so-called *time-harmonic ansatz* for the magnetic field \mathbf{H} . This time-harmonic ansatz is given by

$$\mathbf{H}(t, \mathbf{x}) = \operatorname{Re}[\mathbf{h}_\omega(\mathbf{x})e^{i\omega t}] \quad \text{for all } (t, \mathbf{x}) \in \mathbb{R} \times \mathbb{R}^3, \quad (3.23)$$

where $\mathbf{h}_\omega : \mathbb{R}^3 \rightarrow \mathbb{C}^3$ is a complex vector field. Plugging the time-harmonic ansatz into (3.22), one obtains

$$\begin{cases} \nabla \times \left[\frac{1}{\varepsilon_r} \nabla \times \mathbf{h}_\omega \right] = \frac{\omega^2}{c_0^2} \mu_r \mathbf{h}_\omega & \text{in } \mathbb{R}^3, \\ \nabla \cdot [\mu_r \mathbf{h}_\omega] = 0 & \text{in } \mathbb{R}^3. \end{cases} \quad (3.24)$$

In contrast to the system of equations (3.22), the time-harmonic system of equations (3.24) is posed in \mathbb{R}^3 rather than in $\mathbb{R} \times \mathbb{R}^3$. Notice that (3.24) can be viewed as a constrained eigenvalue problem for the eigenvalue ω^2/c_0^2 and corresponding eigenfunctions \mathbf{h}_ω . The divergence constraint in the second line of (3.24) is natural for the problem in the following sense: Suppose that \mathbf{h}_ω is a complex vector field that solves the eigenvalue equation in the first line of (3.24) for some non-vanishing eigenvalue ω^2/c_0^2 . Then, \mathbf{h}_ω is a curl field, which is why it also satisfies the divergence constraint. For this reason the divergence equation is often omitted in the

literature. We shall keep this equation, nevertheless, as it will play an important role in the analysis and the numerical discretization of the problem.

The constrained eigenvalue problem (3.24) can be further simplified by assuming that the medium, in which the wave propagates, is *non-magnetic*. For such media the relative magnetic permeability is given by

$$\mu_r = 1. \quad (3.25)$$

The system (3.24) hence becomes

$$\begin{cases} \nabla \times \left[\frac{1}{\varepsilon_r} \nabla \times \mathbf{h}_\omega \right] = \frac{\omega^2}{c_0^2} \mathbf{h}_\omega & \text{in } \mathbb{R}^3, \\ \nabla \cdot \mathbf{h}_\omega = 0 & \text{in } \mathbb{R}^3. \end{cases} \quad (3.26)$$

So far, we only considered the magnetic field of an electromagnetic wave propagating in a non-magnetic, non-dispersive, non-absorptive, linear, isotropic dielectric. For completeness, we also introduce the mathematical model for the associated electric field in the following.

When investigating the transition properties of the electric field \mathbf{E} at a smooth surface Σ , which is oriented by the unit normal field $\boldsymbol{\nu}_\Sigma : \Sigma \rightarrow \mathbb{S}^2$, one finds that

$$\boldsymbol{\tau} \cdot (\mathbf{E}^+ - \mathbf{E}^-)|_\Sigma = 0 \quad (3.27)$$

holds in the sense of integrals along smooth curves on Σ for every unit vector field $\boldsymbol{\tau} : \Sigma \rightarrow \mathbb{S}^2$, which is *tangential* to Σ . Recall that the latter notion is made precise by (3.13). As is the case of the magnetic field, we derive from (3.27) that

$$\boldsymbol{\nu}_\Sigma \times (\mathbf{E}^+ - \mathbf{E}^-)|_\Sigma = \mathbf{0},$$

and that

$$\boldsymbol{\nu}_\Sigma \times ((\mathbf{E}^+ - \mathbf{E}^-)|_\Sigma \times \boldsymbol{\nu}_\Sigma) = \mathbf{0}.$$

One can also show that the electric field satisfies

$$\boldsymbol{\nu}_\Sigma \cdot (\varepsilon_r^+ \mathbf{E}^+ - \varepsilon_r^- \mathbf{E}^-)|_\Sigma = 0 \quad (3.28)$$

in the sense of integrals over surfaces on Σ . Clearly, (3.27) and (3.28) are analogues of (3.12) and (3.19), which state that only the tangential component of the electric field is continuous at medium interfaces. The normal component, however, is discontinuous at interfaces, where the relative electric permittivity ε_r changes discontinuously.

One can also derive a second-order, time-harmonic system for the electric field. This is easily done by differentiating Ampère's law in (3.20) with respect to the time variable t , and by then plugging in the identity

$$\dot{\mathbf{H}} = -\frac{1}{\mu_0} \nabla \times \mathbf{E}, \quad (3.29)$$

which follows from Faraday's law. The resulting system reads

$$\begin{cases} \nabla \times \nabla \times \mathbf{E} + \frac{1}{c_0^2} \varepsilon_r \ddot{\mathbf{E}} = 0 & \text{in } \mathbb{R} \times \mathbb{R}^3, \\ \nabla \cdot [\varepsilon_r \mathbf{E}] = 0 & \text{in } \mathbb{R} \times \mathbb{R}^3. \end{cases}, \quad (3.30)$$

With the time-harmonic ansatz

$$\mathbf{E}(t, \mathbf{x}) = \operatorname{Re}[\mathbf{e}_\omega(\mathbf{x})e^{i\omega t}] \quad \text{for all } (t, \mathbf{x}) \in \mathbb{R} \times \mathbb{R}^3, \quad (3.31)$$

where $\mathbf{e}_\omega : \mathbb{R}^3 \rightarrow \mathbb{C}^3$ is a complex vector field, one obtains the system of equations

$$\begin{cases} \nabla \times \nabla \times \mathbf{e}_\omega = \frac{\omega^2}{c_0^2} \varepsilon_r \mathbf{e}_\omega & \text{in } \mathbb{R}^3, \\ \nabla \cdot [\varepsilon_r \mathbf{e}_\omega] = 0 & \text{in } \mathbb{R}^3. \end{cases} \quad (3.32)$$

Clearly, the system of equations (3.32) poses a constrained eigenvalue problem for the eigenvalue ω^2/c_0^2 and corresponding eigenfunctions \mathbf{e}_ω , which is similar to that posed by the system of equations for the magnetic field (3.24).

Finally, we remark that the theory of classical electromagnetism can also be expressed in terms of differential forms on smooth manifolds. In this framework the electric field \mathbf{E} and the magnetic field \mathbf{H} are described by 1-forms which, in essence, represent curve integrals. This alludes to the fact that in practice electric and magnetic fields are always measured with respect to oriented curves in physical space. Recall also that the identities (3.12) and (3.27) should be understood in the sense of integrals along smooth curves at an interface. The electric displacement field \mathbf{D} and the magnetic flux density \mathbf{B} are described by 2-forms, which represent surface integrals. This embodies the fact that the corresponding physical quantities are measured with respect to oriented surfaces in physical space. In view of this, it is not surprising that the identities (3.19) and (3.28) should be understood in the sense of surface integrals, since according to (3.14) and (3.15) the vector fields $\varepsilon_r \mathbf{E}$ and $\mu_r \mathbf{H}$ are proportional to \mathbf{D} and \mathbf{B} , respectively. Expressing electromagnetic waves in terms of differential forms can also be useful, when studying finite element discretizations of Maxwell's equations. By viewing the electric field and the magnetic field as 1-forms, for example, one is able to identify the so-called *Nédélec elements* (cf. [57]) as the class of finite elements, which is best suited for discretizing these fields. It is beyond the scope of this work to make this notion precise. For details, we refer to the works of Arnold [7] and Hiptmair [39].

3.4 Bloch Modes in Periodic Media

The mathematical model which was introduced in the previous section, is valid for time-harmonic wave propagation in any non-magnetic, non-dispersive, non-

absorptive, linear, isotropic, dielectric medium. In this section we consider the special case of wave propagation inside a photonic crystal, whose medium structure, in addition, exhibits some spatial periodicity.

In the following, we assume that the photonic crystal's spatial periodicity is represented by a Bravais lattice $\Lambda \in \mathbb{R}^3$ of rank 3. In Section 3.1 we discussed that the crystal's medium structure can be represented by a Λ -periodic function $\chi : \mathbb{R}^3 \rightarrow \mathbb{N}$, such that every function value $\chi(\mathbf{x})$ represents the material, which is present at the position $\mathbf{x} \in \mathbb{R}^3$. Since each material has its own specific electric permittivity, it follows that the relative electric permittivity field $\varepsilon_r : \mathbb{R}^3 \rightarrow \mathbb{R}$ is also Λ -periodic, i.e.,

$$\varepsilon_r \circ \tau_{\mathbf{a}} = \varepsilon_r \quad \text{for all } \mathbf{a} \in \Lambda.$$

Consequently, (3.26) and (3.32) are systems of partial differential equations with Λ -periodic coefficients for photonic crystals. Such partial differential equations are the subject of the *Floquet–Bloch theory*. This theory is a generalization of the *Floquet theory*. One of the main results of the Floquet theory is Floquet's theorem, which characterizes the fundamental matrices of ordinary, linear differential equations as products of periodic functions and matrix exponentials (cf. [37]). In his work on the motion of electrons in a periodic potential, Bloch showed that the eigenfunctions of Schrödinger-type operators with periodic coefficients are products of periodic functions and plane waves (see [13]). Such functions are now called *Bloch functions*. The findings of Bloch are, in essence, an analog of Floquet's theorem for Schrödinger-type operators. It could be shown that variants of Floquet's theorem apply to a large class of linear partial differential operators, and especially to elliptic ones (see e.g. [49]).

Reconsidering the time-harmonic system of equations for the magnetic field (3.26), the Floquet–Bloch theory suggests that the vector field \mathbf{h}_ω is the product of a Λ -periodic vector field and a plane wave function. More precisely, we make the ansatz

$$\mathbf{h}_\omega = \mathbf{h}_{\omega, \mathbf{k}} e^{i\langle \mathbf{k}, \cdot \rangle}, \quad (3.33)$$

where $\mathbf{h}_{\omega, \mathbf{k}} : \mathbb{R}^3 \rightarrow \mathbb{C}^3$ is a Λ -periodic vector field, and where \mathbf{k} is a vector belonging to the so-called *first Brillouin zone* \mathbb{B} of Λ . As was discussed in Section 2.5, there exists a so-called reciprocal lattice $\widehat{\Lambda} \subset \mathbb{R}^3$ of Λ . The reciprocal lattice $\widehat{\Lambda}$ is a lattice of rank 3, such that

$$\mathbf{a} \cdot \mathbf{b} \in 2\pi\mathbb{Z} \quad \text{for all } \mathbf{a} \in \Lambda, \mathbf{b} \in \widehat{\Lambda}.$$

The first Brillouin zone \mathbb{B} of Λ is defined as the Wigner–Seitz cell of $\widehat{\Lambda}$, i.e.,

$$\mathbb{B} := W_{\widehat{\Lambda}}. \quad (3.34)$$

Since the rank of the reciprocal lattice $\widehat{\Lambda}$ is maximal, the first Brillouin zone is a closed, bounded, convex polyhedron, which is point-symmetric with respect to the origin. The ansatz (3.33) is called a *Bloch ansatz*. The vector \mathbf{k} is called the *quasimomentum vector* of the Bloch function \mathbf{h}_ω . One easily verifies that the Bloch function \mathbf{h}_ω solves the time-harmonic system (3.26), if and only if the function $\mathbf{h}_{\omega,\mathbf{k}}$ satisfies

$$\begin{cases} (\nabla + i\mathbf{k}) \times \left[\frac{1}{\varepsilon_r} (\nabla + i\mathbf{k}) \times \mathbf{h}_{\omega,\mathbf{k}} \right] = \frac{\omega^2}{c_0^2} \mathbf{h}_{\omega,\mathbf{k}} & \text{in } \mathbb{R}^3, \\ (\nabla + i\mathbf{k}) \cdot \mathbf{h}_{\omega,\mathbf{k}} = 0 & \text{in } \mathbb{R}^3. \end{cases} \quad (3.35)$$

The operator $(\nabla + i\mathbf{k}) \times$ is referred to as a *modified curl operator*. The operator $(\nabla + i\mathbf{k}) \cdot$ is called a *modified divergence operator*. The system of equations (3.35) constitutes a family of constrained eigenvalue problems, which is parametrized by the quasimomentum vector \mathbf{k} . Since the vector field $\mathbf{h}_{\omega,\mathbf{k}}$ and the relative electric permittivity function ε_r in (3.35) are Λ -periodic by assumption, it suffices to solve this family of eigenvalue problems in a primitive cell of Λ only. Given such a primitive cell $\Omega \subset \mathbb{R}^3$, we hence consider the family of constrained eigenvalue problems

$$\begin{cases} (\nabla + i\mathbf{k}) \times \left[\frac{1}{\varepsilon_r} (\nabla + i\mathbf{k}) \times \mathbf{h}_{\omega,\mathbf{k}} \right] = \frac{\omega^2}{c_0^2} \mathbf{h}_{\omega,\mathbf{k}} & \text{in } \Omega, \\ (\nabla + i\mathbf{k}) \cdot \mathbf{h}_{\omega,\mathbf{k}} = 0 & \text{in } \Omega, \end{cases} \quad (3.36)$$

for the eigenvalue ω^2/c_0^2 and corresponding eigenfunctions $\mathbf{h}_{\omega,\mathbf{k}}$. This family of constrained eigenvalue problems is indexed by the quasimomentum vector \mathbf{k} , which varies over the first Brillouin zone \mathbb{B} of Λ .

Every eigensolution $(\omega^2/c_0^2, \mathbf{h}_{\omega,\mathbf{k}})$ of (3.36) determines via (3.33) a so-called *Bloch mode* of the photonic crystal with frequency ω and quasimomentum vector \mathbf{k} . In Chapter 4 we will discuss the existence of such Bloch modes. Bloch modes play an essential role in the investigation of the optical properties of photonic crystals. Roughly speaking, an electromagnetic wave with a given frequency ω is able to propagate inside a photonic crystal, if there exists a Bloch mode with the same frequency. In Chapter 4 we will show that for a given quasimomentum vector $\mathbf{k} \in \mathbb{B}$ there exists a countable set of non-negative, real eigenvalues of (3.36) and hence also a countable set of Bloch mode frequencies. The variation of these frequencies with respect to the quasimomentum vector \mathbf{k} is continuous. Hence, the Bloch mode frequencies can be interpreted as continuous functions of the quasimomentum vector. The graphs of these functions are referred to as *photonic bands*. In Chapter 5 we shall study optimization problems related to these photonic bands.

So far, we only considered the magnetic field of an electromagnetic wave propagating in a periodic medium. In the following, we briefly study the associated

electric field. By virtue of the Bloch ansatz

$$\mathbf{e}_\omega = \mathbf{e}_{\omega, \mathbf{k}} e^{i(\mathbf{k}, \cdot)}, \quad (3.37)$$

where $\mathbf{e}_{\omega, \mathbf{k}} : \mathbb{R}^3 \rightarrow \mathbb{C}^3$ is a Λ -periodic vector field, and where $\mathbf{k} \in \mathbb{B}$ is a quasimomentum vector, the time-harmonic system of equations for the electric field (3.32) can be rewritten as

$$\begin{cases} (\nabla + i\mathbf{k}) \times (\nabla + i\mathbf{k}) \times \mathbf{e}_{\omega, \mathbf{k}} = \frac{\omega^2}{c_0^2} \varepsilon_r \mathbf{e}_{\omega, \mathbf{k}} & \text{in } \Omega, \\ (\nabla + i\mathbf{k}) \cdot [\varepsilon_r \mathbf{e}_{\omega, \mathbf{k}}] = 0 & \text{in } \Omega. \end{cases} \quad (3.38)$$

Clearly, the system of equations (3.38) constitutes a family of constrained eigenvalue problems indexed by the quasimomentum vector \mathbf{k} for the eigenvalue ω^2/c_0^2 and corresponding eigenfunctions $\mathbf{e}_{\omega, \mathbf{k}}$. In this work we shall not consider this family of constrained eigenvalue problems for the electric field any further, since for the eigensolutions of (3.38) are determined by those of (3.36), and vice versa. More precisely, one can show that if $(\omega^2/c_0^2, \mathbf{h}_{\omega, \mathbf{k}})$ is an eigensolution of (3.36) with non-vanishing frequency ω for a given quasimomentum vector $\mathbf{k} \in \mathbb{B}$, then $(\omega^2/c_0^2, \mathbf{e}_{\omega, \mathbf{k}})$, where the eigenfunction $\mathbf{e}_{\omega, \mathbf{k}}$ is given by

$$\mathbf{e}_{\omega, \mathbf{k}} = \frac{1}{i\omega\varepsilon_0\varepsilon_r} (\nabla + i\mathbf{k}) \times \mathbf{h}_{\omega, \mathbf{k}},$$

is an eigensolution of (3.38). Conversely, if $(\omega^2/c_0^2, \mathbf{e}_{\omega, \mathbf{k}})$ is an eigensolution of (3.38) with non-vanishing frequency ω for a given quasimomentum vector $\mathbf{k} \in \mathbb{B}$, then $(\omega^2/c_0^2, \mathbf{h}_{\omega, \mathbf{k}})$ with

$$\mathbf{h}_{\omega, \mathbf{k}} = \frac{1}{i\omega\mu_0} (\nabla + i\mathbf{k}) \times \mathbf{e}_{\omega, \mathbf{k}}$$

is an eigensolution of (3.36). It hence suffices to only consider the family of constrained eigenvalue problems for the magnetic field (3.36) in order to determine Bloch mode frequencies for three-dimensional photonic crystals.

3.5 The Two-Dimensional Case

In this section we consider time-harmonic wave propagation in two-dimensional photonic crystals and deduce simplified models for waves with specific polarizations.

In Section 3.1, two-dimensional photonic crystals were characterized as media, which are periodic in two spatial dimensions while being homogeneous in the third

one. By choosing a suitable coordinate system, the crystal function of a two-dimensional photonic crystal can be defined as a function from \mathbb{R}^2 into \mathbb{N} , which is periodic with respect to a Bravais lattice $\Lambda \subset \mathbb{R}^2$ of rank 2. Consequently, the relative electric permittivity ε_r can be identified with a Λ -periodic function mapping from \mathbb{R}^2 into \mathbb{R} . Furthermore, since the electric permittivity only depends on two spatial coordinates, it is reasonable to assume the same for the electric and magnetic field of a wave, which propagates inside a two-dimensional photonic crystal. In this section, we hence assume that the electric and magnetic field are represented by vector fields $\mathbf{E}, \mathbf{H} : \mathbb{R} \times \mathbb{R}^2 \rightarrow \mathbb{R}^3$, which depend on a time variable $t \in \mathbb{R}$ and a spatial variable $\mathbf{x} \in \mathbb{R}^2$. A time-harmonic ansatz for the two vector fields is given by

$$\mathbf{E}(t, \mathbf{x}) = \operatorname{Re}[\mathbf{e}_\omega(\mathbf{x})e^{i\omega t}] \quad \text{for all } (t, \mathbf{x}) \in \mathbb{R} \times \mathbb{R}^2, \quad (3.39)$$

$$\mathbf{H}(t, \mathbf{x}) = \operatorname{Re}[\mathbf{h}_\omega(\mathbf{x})e^{i\omega t}] \quad \text{for all } (t, \mathbf{x}) \in \mathbb{R} \times \mathbb{R}^2 \quad (3.40)$$

where $\mathbf{e}_\omega, \mathbf{h}_\omega : \mathbb{R}^2 \rightarrow \mathbb{C}^3$ are complex vector field, and where $\omega \in \mathbb{R}$ is a phase frequency. With this ansatz, two constrained eigenvalue problems in \mathbb{R}^2 can be established for the vector fields \mathbf{e}_ω and \mathbf{h}_ω in analogy to the time-harmonic systems (3.32) and (3.26).

It turns out the the constrained eigenvalue problems can be simplified significantly, when time-harmonic waves with specific *polarizations* are considered. Here, the term polarization refers to a property of electromagnetic waves, which characterizes the directions of oscillation of the wave's electric and magnetic field. For instance, an electromagnetic wave is called *TM-polarized* or simply a *TM-wave*, if the first two components of its electric field vanish identically. With regard to the time-harmonic ansatz (3.39)–(3.40), an electromagnetic wave is TM-polarized, if and only if the complex vector field \mathbf{e}_ω is of the form

$$\mathbf{e}_\omega = \begin{pmatrix} 0 \\ 0 \\ e_\omega^{\text{TM}} \end{pmatrix}, \quad (3.41)$$

where $e_\omega^{\text{TM}} : \mathbb{R}^2 \rightarrow \mathbb{C}$ is a complex scalar field. According to (3.21) the magnetic field of a time-harmonic TM-wave is determined by the complex vector field

$$\mathbf{h}_\omega^{\text{TM}} := -\frac{1}{\mu_0\omega} \begin{pmatrix} \partial_2 e_\omega^{\text{TM}} \\ -\partial_1 e_\omega^{\text{TM}} \\ 0 \end{pmatrix}, \quad (3.42)$$

where $\partial_i e_\omega^{\text{TM}}$ denotes the partial derivative of the function e_ω^{TM} with respect to the i -th spatial coordinate for $i = 1, 2$. One easily verifies that the time-harmonic

system (3.32) posed in \mathbb{R}^2 reduces for TM-polarized waves to a scalar Laplace-type eigenvalue problem, which reads

$$-\Delta e_\omega^{\text{TM}} = \frac{\omega^2}{c_0^2} \varepsilon_r e_\omega^{\text{TM}} \quad \text{in } \mathbb{R}^2. \quad (3.43)$$

Notice, in particular, that the divergence constraint is automatically fulfilled, since neither ε_r nor e_ω^{TM} depend on the third spatial coordinate.

Another polarization, which affords a simplification of the time-harmonic system (3.26), is the so-called TE-polarization. An electromagnetic wave is called *TE-polarized* or simply a *TE-wave*, if the first two components of its magnetic field vanish identically. Hence, time-harmonic, TE-polarized waves are given by complex vector fields \mathbf{h}_ω , which are of the form

$$\mathbf{h}_\omega = \begin{pmatrix} 0 \\ 0 \\ h_\omega^{\text{TE}} \end{pmatrix}, \quad (3.44)$$

where $h_\omega^{\text{TE}} : \mathbb{R}^2 \rightarrow \mathbb{C}$ is a complex scalar field. According to (3.29), the electric field of a time-harmonic TE-wave is determined by the complex vector field

$$\mathbf{e}_\omega^{\text{TE}} := \frac{1}{\varepsilon_0 \varepsilon_r \omega} \begin{pmatrix} \partial_2 h_\omega^{\text{TE}} \\ -\partial_1 h_\omega^{\text{TE}} \\ 0 \end{pmatrix}. \quad (3.45)$$

For TE-polarized waves the constrained eigenvalue problem (3.26) can be simplified to a second-order, divergence-type eigenvalue problem, which is given by

$$-\nabla \cdot \left[\frac{1}{\varepsilon_r} \nabla h_\omega^{\text{TE}} \right] = \frac{\omega^2}{c_0^2} h_\omega^{\text{TE}} \quad \text{in } \mathbb{R}^2. \quad (3.46)$$

Again, the divergence constraint is fulfilled due to the fact that the scalar field h_ω^{TE} does not depend on the third spatial coordinate. Looking at (3.41), (3.42), (3.44), and (3.45), one realizes that every electromagnetic wave, which propagates inside a two-dimensional photonic crystal, can be represented as the sum of a TM-wave and a TE-wave. Therefore, it is common to consider the simplified eigenvalue problems (3.43) and (3.46) for two-dimensional photonic crystals instead of the more general eigenvalue problems (3.26) and (3.32) posed in \mathbb{R}^2 .

Notice that the eigenvalue problems for the TM- and TE-waves (3.43) and (3.46) are given in terms of partial differential equations with the Λ -periodic coefficients ε_r and $1/\varepsilon_r$, respectively. Therefore, we make a Bloch ansatz

$$e_\omega^{\text{TM}} = e_{\omega, \mathbf{k}}^{\text{TM}} e^{i(\mathbf{k}, \cdot)}, \quad (3.47)$$

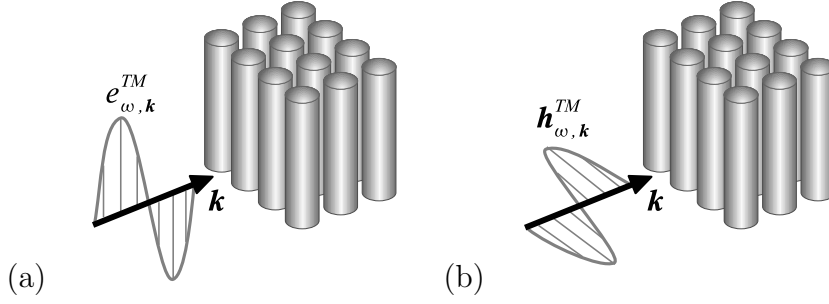


Figure 3.8: Directions of oscillation for the electric field (a) and the magnetic field (b) of a TM-polarized wave in the sense of two-dimensional photonic crystals.

$$h_{\omega}^{\text{TE}} = h_{\omega, \mathbf{k}}^{\text{TE}} e^{i(\mathbf{k}, \cdot)}, \quad (3.48)$$

where $e_{\omega, \mathbf{k}}, h_{\omega, \mathbf{k}} : \mathbb{R}^2 \rightarrow \mathbb{C}$ are Λ -periodic functions, and where \mathbf{k} is a quasimomentum vector belonging to the first Brillouin-zone \mathbb{B} of $\widehat{\Lambda}$. Note that for two-dimensional photonic crystals the first Brillouin \mathbb{B} is a closed, bounded, convex polygon, since the Bravais lattice Λ has maximal rank. Plugging the Bloch ansatz (3.47) into (3.43) yields

$$-(\nabla + i\mathbf{k}) \cdot (\nabla + i\mathbf{k}) e_{\omega, \mathbf{k}}^{\text{TM}} = \frac{\omega^2}{c_0^2} \varepsilon_r e_{\omega, \mathbf{k}}^{\text{TM}} \quad \text{in } \Omega, \quad (3.49)$$

where Ω is an arbitrary primitive cell of Λ . It suffices to solve the eigenvalue problem (3.49) in a primitive cell of Λ , since the functions $e_{\omega, \mathbf{k}}^{\text{TM}}$ and ε_r are Λ -periodic by assumption. Note that (3.49) constitutes a family of eigenvalue problems for TM-waves with the quasimomentum vector \mathbf{k} varying in the first Brillouin zone \mathbb{B} . Similarly, by plugging the Bloch ansatz (3.48) into (3.46), one obtains a family of eigenvalue problems for TE-waves, which reads

$$-(\nabla + i\mathbf{k}) \cdot \left[\frac{1}{\varepsilon_r} (\nabla + i\mathbf{k}) h_{\omega, \mathbf{k}}^{\text{TE}} \right] = \frac{\omega^2}{c_0^2} h_{\omega, \mathbf{k}}^{\text{TE}} \quad \text{in } \Omega. \quad (3.50)$$

As in the case of three-dimensional photonic crystals, a solution $(\omega, e_{\omega, \mathbf{k}}^{\text{TM}})$ of (3.49) determines the electric field of a TM-polarized Bloch mode of the photonic crystal with frequency ω . The corresponding magnetic field $\mathbf{h}_{\omega, \mathbf{k}}^{\text{TM}}$ can be obtained from (3.42) and (3.47). Similarly, a solution $(\omega, e_{\omega, \mathbf{k}}^{\text{TE}})$ determines the magnetic field of a TE-polarized Bloch mode, and the corresponding electric field $\mathbf{e}_{\omega, \mathbf{k}}^{\text{TE}}$ can be obtained from (3.45) and (3.48). An electromagnetic wave with a specific frequency is able to propagate inside a two-dimensional photonic crystal if there exists a TM-polarized Bloch mode or a TE-polarized Bloch mode with the same frequency.

Finally, we mention that TM and TE are abbreviations for *transverse magnetic* and *transverse electric*. The names for the corresponding polarizations originally came from the theory of electromagnetic wave guides, where the propagation of an electromagnetic wave is always assumed in the direction of the third spatial coordinate. The magnetic fields of TM-waves hence oscillate in a plane, which is perpendicular (or transverse) to the direction of propagation. For two-dimensional photonic crystals, however, the term transverse magnetic is rather misleading, since the magnetic fields of TM-waves actually oscillate in the plane of wave propagation as we illustrated in Figure 3.8. The same misconception can arise for TE-waves. Nevertheless, the abbreviations TM and TE are commonly used in the literature on photonic crystals in order to indicate that the polarization of an electromagnetic wave is given by (3.41) or (3.44), respectively.

Chapter 4

Spectral Theory

In Sections 3.3 and 3.4 we developed a mathematical model for time-harmonic wave propagation in three-dimensional photonic crystals. This model leads to a family of constrained eigenvalue problems for the wave's magnetic field, which was given by (3.36). As it turns out the crystal's relative electric permittivity function ε_r enters as a coefficient in each eigenvalue problem. In this chapter, we shall consider the constrained eigenvalue problems from a more general point of view. In Section 4.1 we therefore define a generic family of constrained eigenvalue problems of similar type. In Section 4.2 we introduce suitable function spaces for the corresponding eigenfunctions. Section 4.3 is concerned with establishing a weak formulation of the eigenvalue problems. This weak formulation will be the basis for the spectral theory we present in Section 4.4. In Section 4.5 we introduce a less well-known characterization principle for eigenvalues of an elliptic partial differential operator. This principle will play an important role in the discussion on the existence of solutions of photonic band structure optimization problems in Chapter 5. In Section 4.6 we relate the generic family of eigenvalue problems introduced in Section 4.1 to the band structures of photonic crystals. In particular, we discuss the effect of the non-translational symmetries of a photonic crystal's medium structure on its photonic band structure. In Section 4.7 we briefly comment on how the results of this chapter can be transferred to two-dimensional photonic crystals.

4.1 The Formal Setting

Throughout this chapter, we assume that $\Lambda \subset \mathbb{R}^3$ is a Bravais lattice of rank 3, which represents the translational symmetry of a three-dimensional photonic crystal. By Ω we denote the interior of a bounded, convex, primitive cell of Λ (see Section 2.5). Since the lattice Λ has maximal rank, Ω is an bounded, convex

domain in \mathbb{R}^3 . Furthermore, we denote by \mathbb{B} the first Brillouin zone of Λ . The first Brillouin zone is the Wigner–Seitz cell of the reciprocal lattice $\widehat{\Lambda}$ of Λ and as such a closed, bounded, convex subset of \mathbb{R}^3 . We then consider a family of eigenvalue problems, which is given as follows.

Problem 4.1. *Given a coefficient $\rho : \Omega \rightarrow \mathbb{R}$, which is positive, bounded, and uniformly bounded away from zero almost everywhere on Ω , and a vector $\mathbf{k} \in \mathbb{B}$, find eigenvalues $\lambda \in \mathbb{C}$ and corresponding, Λ -periodic eigenfunctions $\mathbf{u} : \mathbb{R}^3 \rightarrow \mathbb{C}^3$, with $\mathbf{u} \neq \mathbf{0}$, such that*

$$\begin{cases} (\nabla + i\mathbf{k}) \times [\rho(\nabla + i\mathbf{k}) \times \mathbf{u}] = \lambda \mathbf{u} & \text{in } \Omega, \\ (\nabla + i\mathbf{k}) \cdot \mathbf{u} = 0 & \text{in } \Omega. \end{cases} \quad (4.1)$$

Clearly, the constrained eigenvalue problem (4.1) is of the same form as (3.36). In (4.1), the coefficient ρ takes the place of the reciprocal relative electric permittivity $1/\varepsilon_r$. More precisely ρ can be identified with $1/\varepsilon_r|_{\Omega}$. The unknown eigenvalue λ takes the place of the quantity ω^2/c_0^2 and the unknown eigenfunction \mathbf{u} that of the complex vector field $\mathbf{h}_{\omega, \mathbf{k}}$, which determines the corresponding Bloch mode with frequency $\omega = c_0\sqrt{\lambda}$ and quasimomentum vector \mathbf{k} .

In Section 3.1 we assumed that a photonic crystal consists of finitely many different materials. Since every dielectric material possess its own specific electric permittivity, the function ε_r is discontinuous in most cases. Therefore, we only assume ρ as a function, which is bounded almost everywhere on Ω . Because of this, the eigenvalue equation in (4.1) can only be understood in the classical sense, where all functions are assumed to be sufficiently smooth. Therefore, we develop a mathematical framework in the following sections, which allows us to treat (4.1) in a weak sense.

4.2 Sobolev Spaces of Periodic Functions

In this section we introduce suitable function spaces for the analysis of the family of constrained eigenvalue problems stated in Problem 4.1. Recall that the unknown eigenfunctions in (4.1) were required to be Λ -periodic, where Λ is the Bravais lattice describing the translational symmetry of the medium. Since (4.1) is posed in a primitive domain Ω of Λ only, one has to ensure that a non-vanishing vector field, which solves (4.1) for some $\lambda \in \mathbb{C}$, also has a suitable Λ -periodic extension to \mathbb{R}^3 . Furthermore, we require the eigenfunctions to satisfy some weak regularity assumptions, in order to establish a weak formulation of (4.1). Having this in mind, we introduce certain Sobolev spaces of periodic functions in this section. Our aim is to give a consistent characterization of these function spaces in terms of Fourier

series. Moreover, we list a number of important results, which are related to the family of constrained eigenvalue problems stated in Problem 4.1.

In the following, we assume that Ω is a primitive domain of a Bravais lattice $\Lambda \subset \mathbb{R}^3$ of rank 3. We define the Sobolev spaces

$$\mathbf{H}_{\text{per}}(\text{curl}; \Omega) := \{ \mathbf{w}|_{\Omega} \mid \mathbf{w} \in \mathbf{H}_{\text{loc}}(\text{curl}; \mathbb{R}^3), \mathbf{w} \text{ is } \Lambda\text{-periodic} \}, \quad (4.2)$$

$$\mathbf{H}_{\text{per}}(\text{div}; \Omega) := \{ \mathbf{f}|_{\Omega} \mid \mathbf{f} \in \mathbf{H}_{\text{loc}}(\text{div}; \mathbb{R}^3), \mathbf{f} \text{ is } \Lambda\text{-periodic} \}, \quad (4.3)$$

$$H_{\text{per}}^1(\Omega) := \{ q|_{\Omega} \mid q \in H_{\text{loc}}^1(\mathbb{R}^3), q \text{ is } \Lambda\text{-periodic} \}. \quad (4.4)$$

The spaces $\mathbf{H}_{\text{per}}(\text{curl}; \Omega)$, $\mathbf{H}_{\text{per}}(\text{div}; \Omega)$, and $H_{\text{per}}^1(\Omega)$ are Hilbert spaces, which are equipped with the inner products $\langle \cdot, \cdot \rangle_{\text{curl}, \Omega}$, $\langle \cdot, \cdot \rangle_{\text{div}, \Omega}$, and $\langle \cdot, \cdot \rangle_{1, \Omega}$, respectively. Clearly, they are proper linear subspaces of $\mathbf{H}(\text{curl}; \Omega)$, $\mathbf{H}(\text{div}; \Omega)$, and $H^1(\Omega)$. Next, our aim is to derive a suitable characterization of these spaces.

Recall that every Λ -periodic function in $L_{\text{loc}}^2(\mathbb{R}^3)$ is represented by its *Fourier expansion* (see e.g. Section 4.19 in [10]). More precisely,

$$f = \sum_{\mathbf{b} \in \widehat{\Lambda}} \widehat{f}_{\mathbf{b}} e^{i\langle \mathbf{b}, \cdot \rangle} \quad \text{for all } f \in L_{\text{loc}}^2(\mathbb{R}^3)$$

in the sense of $L_{\text{loc}}^2(\Omega)$. The so-called *Fourier coefficients* $\widehat{f}_{\mathbf{b}}$ of f are given by

$$\widehat{f}_{\mathbf{b}} := \frac{1}{\text{meas}(\Omega)} \int_{\Omega} f e^{-i\langle \mathbf{b}, \cdot \rangle} \quad \text{for all } \mathbf{b} \in \widehat{\Lambda}. \quad (4.5)$$

According to *Parseval's identity*, we also have that

$$\|f\|_{\Omega}^2 = \text{meas}(\Omega) \sum_{\mathbf{b} \in \widehat{\Lambda}} |\widehat{f}_{\mathbf{b}}|^2 \quad \text{for all } f \in L_{\text{loc}}^2(\mathbb{R}^3).$$

It should be noted that Fourier coefficients, as defined by (4.5), can be computed for every Λ -periodic function in $L_{\text{loc}}^1(\mathbb{R}^3)$. This is because the restriction of a function in $L_{\text{loc}}^1(\mathbb{R}^3)$ to the primitive domain Ω belongs to $L^1(\Omega)$. Parseval's identity implies that the restriction of a Λ -periodic function $f \in L_{\text{loc}}^1(\mathbb{R}^3)$ to Ω belongs to $L^2(\Omega)$, if and only if its Fourier coefficients satisfy

$$\sum_{\mathbf{b} \in \widehat{\Lambda}} |\widehat{f}_{\mathbf{b}}|^2 < \infty. \quad (4.6)$$

If that is the case, the function f obviously belongs to $L_{\text{loc}}^2(\mathbb{R}^3)$ due to the Λ -periodicity. It follows, that a Λ -periodic function $f : \mathbb{R}^3 \rightarrow \mathbb{C}$ belongs to $L_{\text{loc}}^2(\mathbb{R}^3)$, if and only if its Fourier coefficients $\widehat{f}_{\mathbf{b}}$ are well-defined, and if they satisfy (4.6).

Now suppose that $g : \Omega \rightarrow \mathbb{C}$ is a some given function. The question is, whether or not g can be understood as the restriction to Ω of a Λ -periodic function $L^2_{\text{loc}}(\Omega)$. If g belongs to $L^1(\Omega)$, one can compute Fourier coefficients $\widehat{g}_{\mathbf{b}}$ of g according to (4.5). By Parseval's identity, the function g even belongs to $L^2(\Omega)$, if and only if

$$\sum_{\mathbf{b} \in \widehat{\Lambda}} |\widehat{g}_{\mathbf{b}}|^2 < \infty. \quad (4.7)$$

Provided that the above estimate holds, one can define a function $\widetilde{g} \in L^2_{\text{loc}}(\mathbb{R}^3)$ by

$$\widetilde{g} := \sum_{\mathbf{b} \in \widehat{\Lambda}} \widehat{g}_{\mathbf{b}} e^{i\langle \mathbf{b}, \cdot \rangle}.$$

Obviously, the function \widetilde{g} is Λ -periodic, and we have that the restriction of \widetilde{g} to the primitive domain Ω coincides with g . Hence, a function $g : \Omega \rightarrow \mathbb{C}$ can be understood as the restriction to Ω of a Λ -periodic function in $L^2_{\text{loc}}(\Omega)$, if and only if its Fourier coefficients $\widehat{g}_{\mathbf{b}}$ are well-defined, and if they satisfy (4.7).

The above discussion shows that the restrictions of Λ -periodic functions in $L^2_{\text{loc}}(\mathbb{R}^3)$ to a primitive domain Ω can be characterized by their Fourier coefficients. In the following, this principle is used to characterize the Sobolev spaces Sobolev spaces $\mathbf{H}_{\text{per}}(\text{curl}; \Omega)$, $\mathbf{H}_{\text{per}}(\text{div}; \Omega)$, and $H^1_{\text{per}}(\Omega)$, which were defined by (4.2)–(4.4).

Given some Λ -periodic functions $\mathbf{w} \in \mathbf{H}_{\text{loc}}(\text{curl}; \mathbb{R}^3)$, $\mathbf{f} \in \mathbf{H}_{\text{loc}}(\text{div}; \mathbb{R}^3)$, and $q \in H^1_{\text{loc}}(\mathbb{R}^3)$, we have that $\nabla \times \mathbf{w}$, $\nabla \cdot \mathbf{f}$, and ∇q belong to $L^2_{\text{loc}}(\mathbb{R}^3)^3$, $L^2_{\text{loc}}(\mathbb{R}^3)$, and $L^2_{\text{loc}}(\mathbb{R}^3)^3$, respectively. Furthermore, we have that these functions are Λ -periodic. Hence, they are represented by their Fourier expansions. More precisely, we have that the identities

$$\nabla \times \mathbf{w} = \sum_{\mathbf{b} \in \widehat{\Lambda}} i\mathbf{b} \times \widehat{\mathbf{w}}_{\mathbf{b}} e^{i\langle \mathbf{b}, \cdot \rangle}, \quad (4.8)$$

$$\nabla \cdot \mathbf{f} = \sum_{\mathbf{b} \in \widehat{\Lambda}} i\mathbf{b} \cdot \widehat{\mathbf{f}}_{\mathbf{b}} e^{i\langle \mathbf{b}, \cdot \rangle}, \quad (4.9)$$

$$\nabla q = \sum_{\mathbf{b} \in \widehat{\Lambda}} i\mathbf{b} \widehat{q}_{\mathbf{b}} e^{i\langle \mathbf{b}, \cdot \rangle} \quad (4.10)$$

hold in the sense of $L^2_{\text{loc}}(\mathbb{R}^3)^3$, $L^2_{\text{loc}}(\mathbb{R}^3)$, and $L^2_{\text{loc}}(\mathbb{R}^3)^3$, respectively. By Parseval's identity we also have that

$$\|\nabla \times \mathbf{w}\|_{\Omega}^2 = \text{meas}(\Omega) \sum_{\mathbf{b} \in \widehat{\Lambda}} |\mathbf{b} \times \widehat{\mathbf{w}}_{\mathbf{b}}|^2, \quad (4.11)$$

$$\|\nabla \cdot \mathbf{f}\|_{\Omega}^2 = \text{meas}(\Omega) \sum_{\mathbf{b} \in \widehat{\Lambda}} |\mathbf{b} \cdot \widehat{\mathbf{f}}_{\mathbf{b}}|^2, \quad (4.12)$$

$$\|\nabla q\|_{\Omega}^2 = \text{meas}(\Omega) \sum_{\mathbf{b} \in \widehat{\Lambda}} |\mathbf{b}|^2 |\widehat{q}_{\mathbf{b}}|^2. \quad (4.13)$$

Since the restrictions of $\nabla \times \mathbf{w}$, $\nabla \cdot \mathbf{f}$, and ∇q to the primitive domain Ω belong to $L^2(\Omega)^3$, $L^2(\Omega)$, and $L^2(\Omega)^3$, it follows that

$$\sum_{\mathbf{b} \in \widehat{\Lambda}} |\mathbf{b} \times \widehat{\mathbf{w}}_{\mathbf{b}}|^2 < \infty, \quad (4.14)$$

$$\sum_{\mathbf{b} \in \widehat{\Lambda}} |\mathbf{b} \cdot \widehat{\mathbf{f}}_{\mathbf{b}}|^2 < \infty, \quad (4.15)$$

$$\sum_{\mathbf{b} \in \widehat{\Lambda}} |\mathbf{b}|^2 |\widehat{q}_{\mathbf{b}}|^2 < \infty. \quad (4.16)$$

Conversely, suppose that \mathbf{w} , \mathbf{f} and q are functions in $L^2(\Omega)^3$, in $L^2(\Omega)^3$, and in $L^2(\Omega)$, whose Fourier coefficients satisfy (4.14)–(4.16). Then, one can extend these functions Λ -periodically to \mathbb{R}^3 using their Fourier series, and thus obtain functions that belong to $\mathbf{H}_{\text{loc}}(\text{curl}; \mathbb{R}^3)$, $\mathbf{H}_{\text{loc}}(\text{div}; \mathbb{R}^3)$, and $H_{\text{loc}}^1(\mathbb{R}^3)$. Hence, the Sobolev spaces $\mathbf{H}_{\text{per}}(\text{curl}; \Omega)$, $\mathbf{H}_{\text{per}}(\text{div}; \Omega)$, and $H_{\text{per}}^1(\Omega)$ can be characterized as

$$\begin{aligned} \mathbf{H}_{\text{per}}(\text{curl}; \Omega) &= \left\{ \mathbf{w} \in L^2(\Omega)^3 \left| \sum_{\mathbf{b} \in \widehat{\Lambda}} (|\widehat{\mathbf{w}}_{\mathbf{b}}|^2 + |\mathbf{b} \times \widehat{\mathbf{w}}_{\mathbf{b}}|^2) < \infty \right. \right\}, \\ \mathbf{H}_{\text{per}}(\text{div}; \Omega) &= \left\{ \mathbf{f} \in L^2(\Omega)^3 \left| \sum_{\mathbf{b} \in \widehat{\Lambda}} (|\widehat{\mathbf{f}}_{\mathbf{b}}|^2 + |\mathbf{b} \cdot \widehat{\mathbf{f}}_{\mathbf{b}}|^2) < \infty \right. \right\}, \\ H_{\text{per}}^1(\Omega) &= \left\{ q \in L^2(\Omega) \left| \sum_{\mathbf{b} \in \widehat{\Lambda}} (1 + |\mathbf{b}|^2) |\widehat{q}_{\mathbf{b}}|^2 < \infty \right. \right\}. \end{aligned}$$

The fact that all functions, which belong to the above Sobolev spaces, are represented by their Fourier expansions greatly simplifies the proofs of the results, which are given below. As we shall see, some of these proofs can be accomplished by relatively simple computations involving Fourier coefficients. Consider, for example, the proof of the following lemma.

Lemma 4.2. *For every vector field $\boldsymbol{\psi} \in H_{\text{per}}^1(\Omega)^3$ the following identity holds,*

$$\|\boldsymbol{\psi}\|_{1,\Omega}^2 = \|\boldsymbol{\psi}\|_{\Omega}^2 + \|\nabla \cdot \boldsymbol{\psi}\|_{\Omega}^2 + \|\nabla \times \boldsymbol{\psi}\|_{\Omega}^2.$$

Proof. Let $\boldsymbol{\psi} = (\psi^{(1)}, \psi^{(2)}, \psi^{(3)}) \in H_{\text{per}}^1(\Omega)^3$ be an arbitrarily chosen vector field. Then, we obtain by (4.13) that

$$\|\boldsymbol{\psi}\|_{1,\Omega}^2 = \|\boldsymbol{\psi}\|_{\Omega}^2 + \|\nabla \psi^{(1)}\|_{\Omega}^2 + \|\nabla \psi^{(2)}\|_{\Omega}^2 + \|\nabla \psi^{(3)}\|_{\Omega}^2$$

$$\begin{aligned}
&= \|\boldsymbol{\psi}\|_{\Omega}^2 + \text{meas}(\Omega) \sum_{\mathbf{b} \in \widehat{\Lambda}} |\mathbf{b}|^2 \left(|\widehat{\psi}^{(1)}_{\mathbf{b}}|^2 + |\widehat{\psi}^{(2)}_{\mathbf{b}}|^2 + |\widehat{\psi}^{(3)}_{\mathbf{b}}|^2 \right) \\
&= \|\boldsymbol{\psi}\|_{\Omega}^2 + \text{meas}(\Omega) \sum_{\mathbf{b} \in \widehat{\Lambda}} |\mathbf{b}|^2 \left(|\widehat{\psi}^{(1)}_{\mathbf{b}}|^2 + |\widehat{\psi}^{(2)}_{\mathbf{b}}|^2 + |\widehat{\psi}^{(3)}_{\mathbf{b}}|^2 \right) \\
&= \|\boldsymbol{\psi}\|_{\Omega}^2 + \text{meas}(\Omega) \sum_{\mathbf{b} \in \widehat{\Lambda}} |\mathbf{b}|^2 |\widehat{\boldsymbol{\psi}}_{\mathbf{b}}|^2.
\end{aligned}$$

Using the vector norm identity (2.4) on page 15, one obtains that $|\mathbf{b}|^2 |\widehat{\boldsymbol{\psi}}_{\mathbf{b}}|^2 = |\mathbf{b} \cdot \widehat{\boldsymbol{\psi}}_{\mathbf{b}}|^2 + |\mathbf{b} \times \widehat{\boldsymbol{\psi}}_{\mathbf{b}}|^2$ holds for every $\mathbf{b} \in \widehat{\Lambda}$. Hence, by (4.11) and (4.12), we have that

$$\begin{aligned}
\|\boldsymbol{\psi}\|_{1,\Omega}^2 &= \|\boldsymbol{\psi}\|_{\Omega}^2 + \text{meas}(\Omega) \sum_{\mathbf{b} \in \widehat{\Lambda}} \left(|\mathbf{b} \cdot \widehat{\boldsymbol{\psi}}_{\mathbf{b}}|^2 + |\mathbf{b} \times \widehat{\boldsymbol{\psi}}_{\mathbf{b}}|^2 \right) \\
&= \|\boldsymbol{\psi}\|_{\Omega}^2 + \|\nabla \cdot \boldsymbol{\psi}\|_{\Omega}^2 + \|\nabla \times \boldsymbol{\psi}\|_{\Omega}^2,
\end{aligned}$$

which completes the proof. \square

It should be noted that Lemma 4.2 implies that

$$H_{\text{per}}^1(\Omega)^3 = \mathbf{H}_{\text{per}}(\text{curl}; \Omega) \cap \mathbf{H}_{\text{per}}(\text{div}; \Omega).$$

Next, we investigate some properties of the modified curl, divergence, and gradient operators $(\nabla + \mathbf{i}\mathbf{k})\times$, $(\nabla + \mathbf{i}\mathbf{k})\cdot$ and $(\nabla + \mathbf{i}\mathbf{k})$, where \mathbf{k} is a vector belonging to the first Brillouin zone \mathbb{B} of Λ . As a first result, we state the following lemma, which establishes important norm estimates.

Lemma 4.3. *For every vector $\mathbf{k} \in \mathbb{B}$ we have the norm estimates*

$$\begin{aligned}
\|(\nabla + \mathbf{i}\mathbf{k}) \times \mathbf{w}\|_{\Omega} &\leq \sqrt{\beta_0(\mathbf{k})} \|\mathbf{w}\|_{\text{curl},\Omega} && \text{for all } \mathbf{w} \in \mathbf{H}_{\text{per}}(\text{curl}; \Omega), \\
\|(\nabla + \mathbf{i}\mathbf{k}) \cdot \mathbf{f}\|_{\Omega} &\leq \sqrt{\beta_0(\mathbf{k})} \|\mathbf{f}\|_{\text{div},\Omega} && \text{for all } \mathbf{f} \in \mathbf{H}_{\text{per}}(\text{div}; \Omega), \\
\|(\nabla + \mathbf{i}\mathbf{k})q\|_{\Omega} &\leq \sqrt{\beta_0(\mathbf{k})} \|q\|_{1,\Omega} && \text{for all } q \in H_{\text{per}}^1(\Omega),
\end{aligned}$$

where the function $\beta_0 : \mathbb{B} \rightarrow \mathbb{R}_{>0}$ is given by

$$\beta_0(\mathbf{k}) := 2 \max\{1, |\mathbf{k}|^2\} \quad \text{for all } \mathbf{k} \in \mathbb{B}.$$

Proof. Given an arbitrary function $\mathbf{w} \in \mathbf{H}_{\text{per}}(\text{curl}, \Omega)$, one obtains

$$\|(\nabla + \mathbf{i}\mathbf{k}) \times \mathbf{w}\|_{\Omega}^2 = \int_{\Omega} |(\nabla + \mathbf{i}\mathbf{k}) \times \mathbf{w}|^2 \leq \int_{\Omega} (|\nabla \times \mathbf{w}| + |\mathbf{k} \times \mathbf{w}|)^2$$

$$\begin{aligned}
&\leq 2 \left(\int_{\Omega} |\nabla \times \mathbf{w}|^2 + \int_{\Omega} |\mathbf{k} \times \mathbf{w}|^2 \right) \\
&= 2 \left(\int_{\Omega} |\nabla \times \mathbf{w}|^2 + |\mathbf{k}|^2 \int_{\Omega} |\mathbf{w}|^2 \right) \\
&\leq 2 \max\{1, |\mathbf{k}|^2\} \|\mathbf{w}\|_{\text{curl}, \Omega}^2,
\end{aligned}$$

which proves the first norm estimate. The remaining norm estimates can be verified analogously. \square

Lemma 4.3 in particular implies that the modified curl, divergence, and gradient operators are continuous, linear operators from $\mathbf{H}_{\text{per}}(\text{curl}, \Omega)$, $\mathbf{H}_{\text{per}}(\text{div}, \Omega)$, and $H^1(\Omega)$ into $L^2(\Omega)^3$, $L^2(\Omega)^3$, and $L^2(\Omega)$, respectively. Hence, given some arbitrary functions $\mathbf{w} \in \mathbf{H}_{\text{per}}(\text{curl}; \Omega)$, $\mathbf{f} \in \mathbf{H}_{\text{per}}(\text{div}; \Omega)$, and $q \in H_{\text{per}}^1(\Omega)$, we have that the functions $(\nabla + i\mathbf{k}) \times \mathbf{w}$, $(\nabla + i\mathbf{k}) \cdot \mathbf{f}$, and $(\nabla + i\mathbf{k})q$ are represented by their Fourier expansions, i.e., we have that

$$(\nabla + i\mathbf{k}) \times \mathbf{w} = \sum_{\mathbf{b} \in \widehat{\Lambda}} i(\mathbf{b} + \mathbf{k}) \times \widehat{\mathbf{w}}_{\mathbf{b}} e^{i(\mathbf{b}, \cdot)}, \quad (4.17)$$

$$(\nabla + i\mathbf{k}) \cdot \mathbf{f} = \sum_{\mathbf{b} \in \widehat{\Lambda}} i(\mathbf{b} + \mathbf{k}) \cdot \widehat{\mathbf{f}}_{\mathbf{b}} e^{i(\mathbf{b}, \cdot)}, \quad (4.18)$$

$$(\nabla + i\mathbf{k})q = \sum_{\mathbf{b} \in \widehat{\Lambda}} i(\mathbf{b} + i\mathbf{k}) \widehat{q}_{\mathbf{b}} e^{i(\mathbf{b}, \cdot)} \quad (4.19)$$

in the sense of $L^2(\Omega)^3$, $L^2(\Omega)$, and $L^2(\Omega)^3$. Using Parseval's identity again, we also find that

$$\|(\nabla + i\mathbf{k}) \times \mathbf{w}\|_{\Omega}^2 = \sum_{\mathbf{b} \in \widehat{\Lambda}} |(\mathbf{b} + \mathbf{k}) \times \widehat{\mathbf{w}}_{\mathbf{b}}|^2, \quad (4.20)$$

$$\|(\nabla + i\mathbf{k}) \cdot \mathbf{f}\|_{\Omega}^2 = \sum_{\mathbf{b} \in \widehat{\Lambda}} |(\mathbf{b} + \mathbf{k}) \cdot \widehat{\mathbf{f}}_{\mathbf{b}}|^2, \quad (4.21)$$

$$\|(\nabla + i\mathbf{k})q\|_{\Omega}^2 = \sum_{\mathbf{b} \in \widehat{\Lambda}} |(\mathbf{b} + \mathbf{k}) \widehat{q}_{\mathbf{b}}|^2. \quad (4.22)$$

The following lemma provides *Gårding*-type inequalities for the modified curl, divergence and gradient operators. We remark that the proof of this lemma relies on the fact that the first Brillouin zone \mathbb{B} is defined as the Wigner–Seitz cell of the reciprocal lattice $\widehat{\Lambda}$. This implies that \mathbb{B} is point symmetric with respect to the origin (see Section 2.5). Therefore, we have that

$$|\mathbf{k}| \leq \frac{1}{2} \text{diam}(\mathbb{B}) \quad \text{for all } \mathbf{k} \in \mathbb{B}, \quad (4.23)$$

where $\text{diam}(\mathbb{B})$ denotes the diameter of \mathbb{B} .

Lemma 4.4. *For every vector $\mathbf{k} \in \mathbb{B}$ we have the inequalities*

$$\begin{aligned} \|(\nabla + i\mathbf{k}) \times \mathbf{w}\|_{\Omega}^2 &\geq \alpha_0(\mathbf{k}) \|\mathbf{w}\|_{\text{curl},\Omega}^2 - \kappa_0(\mathbf{k}) \|\mathbf{w}\|_{\Omega}^2 && \text{for all } \mathbf{w} \in \mathbf{H}_{\text{per}}(\text{curl}; \Omega), \\ \|(\nabla + i\mathbf{k}) \cdot \mathbf{f}\|_{\Omega}^2 &\geq \alpha_0(\mathbf{k}) \|\mathbf{f}\|_{\text{div},\Omega}^2 - \kappa_0(\mathbf{k}) \|\mathbf{f}\|_{\Omega}^2 && \text{for all } \mathbf{f} \in \mathbf{H}_{\text{per}}(\text{div}; \Omega), \\ \|(\nabla + i\mathbf{k})q\|_{\Omega}^2 &\geq \alpha_0(\mathbf{k}) \|q\|_{1,\Omega}^2 - \kappa_0(\mathbf{k}) \|q\|_{\Omega}^2 && \text{for all } q \in H_{\text{per}}^1(\Omega), \end{aligned}$$

where the functions $\alpha_0 : \mathbb{B} \rightarrow \mathbb{R}_{>0}$ and $\kappa_0 : \mathbb{B} \rightarrow \mathbb{R}_{>0}$ are given by

$$\begin{aligned} \alpha_0(\mathbf{k}) &:= \left(1 - \frac{|\mathbf{k}|}{\text{diam}(\mathbb{B})}\right) && \text{for all } \mathbf{k} \in \mathbb{B}. \\ \kappa_0(\mathbf{k}) &:= \alpha_0(\mathbf{k}) + \text{diam}(\mathbb{B})|\mathbf{k}| - |\mathbf{k}|^2 && \text{for all } \mathbf{k} \in \mathbb{B}. \end{aligned}$$

Proof. Let $\mathbf{w} \in \mathbf{H}_{\text{per}}(\text{curl}; \Omega)$ be chosen arbitrarily. Then, we have

$$\begin{aligned} \|(\nabla + i\mathbf{k}) \times \mathbf{w}\|_{\Omega}^2 &= \int_{\Omega} |(\nabla + i\mathbf{k}) \times \mathbf{u}|^2 \\ &= \int_{\Omega} \left(|\nabla \times \mathbf{u}|^2 + 2 \text{Re}[\overline{\nabla \times \mathbf{u}} \cdot (i\mathbf{k} \times \mathbf{u})] + |\mathbf{k} \times \mathbf{u}|^2 \right) \\ &\geq \int_{\Omega} \left(|\nabla \times \mathbf{u}|^2 - 2|\nabla \times \mathbf{u}||\mathbf{k} \times \mathbf{u}| + |\mathbf{k} \times \mathbf{u}|^2 \right) \end{aligned}$$

By Young's inequality, we have that

$$|\nabla \times \mathbf{u}|_{\Omega} |\mathbf{k} \times \mathbf{u}|_{\Omega} \leq \frac{1}{2} \left(\frac{1}{\theta} |\nabla \times \mathbf{u}|_{\Omega}^2 + \theta |\mathbf{k} \times \mathbf{u}|_{\Omega}^2 \right) \quad \text{for all } \theta > 0.$$

Letting $\theta := \text{diam}(\mathbb{B})/|\mathbf{k}|$, we have that $\theta \geq 2$ according to (4.23). Using the vector norm estimate (2.4) on page 15, we hence get

$$\begin{aligned} \|(\nabla + i\mathbf{k}) \times \mathbf{w}\|_{\Omega}^2 &\geq \left(1 - \frac{1}{\theta}\right) \int_{\Omega} |\nabla \times \mathbf{u}|^2 - (\theta - 1) \int_{\Omega} |\mathbf{k} \times \mathbf{u}|^2 \\ &\geq \left(1 - \frac{1}{\theta}\right) \int_{\Omega} |\nabla \times \mathbf{u}|^2 - (\theta - 1) \int_{\Omega} |\mathbf{k}|^2 |\mathbf{u}|^2 \\ &= \left(1 - \frac{|\mathbf{k}|}{\text{diam}(\mathbb{B})}\right) \|\nabla \times \mathbf{u}\|_{\Omega}^2 - (\text{diam}(\mathbb{B})|\mathbf{k}| - |\mathbf{k}|^2) \|\mathbf{u}\|_{\Omega}^2 \end{aligned}$$

which establishes the inequality stated in (a). The remaining inequalities can be established through a similar line of argument. \square

Note that the functions α_0 and κ_0 in Lemma 4.4 are indeed positive everywhere on \mathbb{B} . This is due to the estimate given in (4.23), which implies that

$$\begin{aligned}\alpha_0(\mathbf{k}) &\geq \frac{1}{2} && \text{for all } \mathbf{k} \in \mathbb{B}, \\ \kappa_0(\mathbf{k}) &\geq \frac{1}{2} && \text{for all } \mathbf{k} \in \mathbb{B}.\end{aligned}$$

We remark that the parameter θ in the proof of Lemma 4.4 should be chosen such that $\alpha_0(\mathbf{0}) = \kappa_0(\mathbf{0}) = 1$. With this, the asserted inequalities become equalities for the case $\mathbf{k} = \mathbf{0}$ as one would expect.

In order to establish a weak formulation of the eigenvalue problem stated in Problem 4.1, it is desirable to have *Green-type* identities for the modified curl, divergence and gradient operators. Such identities are established in the following theorem.

Theorem 4.5. *The following assertions hold for every vector $\mathbf{k} \in \mathbb{B}$.*

(a) *For every $\mathbf{v}, \mathbf{w} \in \mathbf{H}_{\text{per}}(\text{curl}; \Omega)$ we have*

$$\int_{\Omega} \overline{(\nabla + i\mathbf{k}) \times \mathbf{w} \cdot \mathbf{v}} = \int_{\Omega} \overline{\mathbf{w}} \cdot (\nabla + i\mathbf{k}) \times \mathbf{v}.$$

(b) *For every $\mathbf{f} \in \mathbf{H}_{\text{per}}(\text{div}; \Omega)$ and $q \in H_{\text{per}}^1(\Omega)$ we have*

$$\int_{\Omega} \overline{(\nabla + i\mathbf{k}) \cdot \mathbf{f}} q = - \int_{\Omega} \overline{\mathbf{f}} \cdot (\nabla + i\mathbf{k}) q.$$

Proof. To prove the identities, we expand all functions according (4.17)–(4.19). Because of (4.20)–(4.22) integration and summation can be interchanged. Since

$$\int_{\Omega} e^{i\langle \mathbf{b} - \tilde{\mathbf{b}}, \cdot \rangle} = \begin{cases} \text{meas}(\Omega) & \text{if } \mathbf{b} = \tilde{\mathbf{b}}, \\ 0 & \text{else} \end{cases} \quad \text{for all } \mathbf{b}, \tilde{\mathbf{b}} \in \hat{\Lambda},$$

we obtain

$$\begin{aligned}\int_{\Omega} \overline{(\nabla + i\mathbf{k}) \times \mathbf{w} \cdot \mathbf{v}} &= - \text{meas}(\Omega) i \sum_{\mathbf{b} \in \hat{\Lambda}} ((\mathbf{b} + \mathbf{k}) \times \overline{\hat{\mathbf{w}}_{\mathbf{b}}}) \cdot \hat{\mathbf{v}}_{\mathbf{b}}, \\ \int_{\Omega} \overline{\mathbf{w}} \cdot (\nabla + i\mathbf{k}) \times \mathbf{v} &= \text{meas}(\Omega) i \sum_{\mathbf{b} \in \hat{\Lambda}} \overline{\hat{\mathbf{w}}_{\mathbf{b}}} \cdot ((\mathbf{b} + \mathbf{k}) \times \hat{\mathbf{v}}_{\mathbf{b}}), \\ \int_{\Omega} \overline{(\nabla + i\mathbf{k}) \cdot \mathbf{f}} q &= - \text{meas}(\Omega) i \sum_{\mathbf{b} \in \hat{\Lambda}} (\mathbf{b} + \mathbf{k}) \cdot \overline{\hat{\mathbf{f}}_{\mathbf{b}}} \hat{q}_{\mathbf{b}},\end{aligned}$$

$$-\int_{\Omega} \overline{\mathbf{f}} \cdot (\nabla + i\mathbf{k})q = \text{meas}(\Omega) i \sum_{\mathbf{b} \in \widehat{\Lambda}} \overline{\widehat{\mathbf{f}}_{\mathbf{b}}} \cdot (\mathbf{b} + \mathbf{k}) \widehat{q}_{\mathbf{b}}.$$

The assertions then follow from the vector identity (2.2) and the fact that $\mathbf{x} \times \mathbf{y} = -\mathbf{y} \times \mathbf{x}$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{C}^3$. \square

It should be noted that in Theorem 4.5 integrals over the boundary $\partial\Omega$ of the fundamental cell Ω do not appear, neither in the integral identities nor in the proof. Again, this is due to the fact that all functions are represented by corresponding Fourier expansions, so that the integral identities reduce to simple vector identities for the Fourier coefficients. Using the same principle, it is easy to verify the following lemma, which extends some well-known results about certain compositions of curl, divergence and gradient operators to their modified counterparts.

Lemma 4.6. *The following assertions hold for every vector $\mathbf{k} \in \mathbb{B}$.*

(a) *For every function $q \in H_{\text{per}}^1(\Omega)$ we have $(\nabla + i\mathbf{k})q \in \mathbf{H}_{\text{per}}(\text{curl}; \Omega)$ and*

$$(\nabla + i\mathbf{k}) \times (\nabla + i\mathbf{k})q = \mathbf{0}.$$

(b) *For every function $\mathbf{w} \in \mathbf{H}_{\text{per}}(\text{curl}; \Omega)$ we have $(\nabla + i\mathbf{k}) \times \mathbf{w} \in \mathbf{H}_{\text{per}}(\text{div}; \Omega)$ and*

$$(\nabla + i\mathbf{k}) \cdot (\nabla + i\mathbf{k}) \times \mathbf{w} = 0.$$

Proof. Let $q \in H_{\text{per}}^1(\Omega)$ and $\mathbf{w} \in \mathbf{H}_{\text{per}}(\text{curl}; \Omega)$ be chosen arbitrarily. Then, by (4.17) and (4.19) the Fourier coefficients of the functions $\mathbf{u} := (\nabla + i\mathbf{k})q$ and $\mathbf{f} := (\nabla + i\mathbf{k}) \times \mathbf{w}$ are given by $\widehat{\mathbf{u}}_{\mathbf{b}} = i(\mathbf{b} + \mathbf{k})\widehat{q}_{\mathbf{b}}$ and $\widehat{\mathbf{f}}_{\mathbf{b}} = i(\mathbf{b} + \mathbf{k}) \times \widehat{\mathbf{w}}_{\mathbf{b}}$ for all $\mathbf{b} \in \widehat{\Lambda}$, where $\widehat{q}_{\mathbf{b}}$ and $\widehat{\mathbf{w}}_{\mathbf{b}}$ denote the respective Fourier coefficients of q and \mathbf{w} . The functions \mathbf{u} and \mathbf{f} belong to $\mathbf{H}_{\text{per}}(\text{curl}; \Omega)$ and $\mathbf{H}_{\text{per}}(\text{div}; \Omega)$, respectively, since

$$\begin{aligned} \sum_{\mathbf{b} \in \widehat{\Lambda}} (|\widehat{\mathbf{u}}_{\mathbf{b}}|^2 + |\mathbf{b} \times \widehat{\mathbf{u}}_{\mathbf{b}}|^2) &= \sum_{\mathbf{b} \in \widehat{\Lambda}} (|(\mathbf{b} + \mathbf{k})\widehat{q}_{\mathbf{b}}|^2 + |\mathbf{b} \times \mathbf{k}\widehat{q}_{\mathbf{b}}|^2) \\ &\leq \sum_{\mathbf{b} \in \widehat{\Lambda}} |(\mathbf{b} + \mathbf{k})\widehat{q}_{\mathbf{b}}|^2 + |\mathbf{k}|^2 \sum_{\mathbf{b} \in \widehat{\Lambda}} |\mathbf{b}|^2 |\widehat{q}_{\mathbf{b}}|^2 < \infty, \end{aligned}$$

and

$$\begin{aligned} \sum_{\mathbf{b} \in \widehat{\Lambda}} (|\widehat{\mathbf{f}}_{\mathbf{b}}|^2 + |\mathbf{b} \cdot \widehat{\mathbf{f}}_{\mathbf{b}}|^2) &= \sum_{\mathbf{b} \in \widehat{\Lambda}} (|(\mathbf{b} + \mathbf{k}) \times \widehat{\mathbf{w}}_{\mathbf{b}}|^2 + |\mathbf{b} \cdot ((\mathbf{b} + \mathbf{k}) \times \widehat{\mathbf{w}}_{\mathbf{b}})|^2) \\ &= \sum_{\mathbf{b} \in \widehat{\Lambda}} (|(\mathbf{b} + \mathbf{k}) \times \widehat{\mathbf{w}}_{\mathbf{b}}|^2 + |\mathbf{k} \cdot (\mathbf{b} \times \widehat{\mathbf{w}}_{\mathbf{b}})|^2) \end{aligned}$$

$$\leq \sum_{\mathbf{b} \in \widehat{\Lambda}} |(\mathbf{b} + \mathbf{k}) \times \widehat{\mathbf{w}}_{\mathbf{b}}|^2 + |\mathbf{k}|^2 \sum_{\mathbf{b} \in \widehat{\Lambda}} |\mathbf{b} \times \widehat{\mathbf{w}}_{\mathbf{b}}|^2 < \infty.$$

To obtain the above inequality we used the same vector identities as in the proof of Theorem 4.8. By (4.17) and (4.18), the Fourier coefficients of the function $\mathbf{v} := (\nabla + \mathbf{i}\mathbf{k}) \times \mathbf{u}$ and $g := (\nabla + \mathbf{i}\mathbf{k}) \cdot \mathbf{f}$ are given by $\widehat{\mathbf{v}}_{\mathbf{b}} = \mathbf{i}(\mathbf{b} + \mathbf{k}) \times \widehat{\mathbf{u}}_{\mathbf{b}} = -(\mathbf{b} + \mathbf{k}) \times (\mathbf{b} + \mathbf{k})\widehat{q}_{\mathbf{b}} = \mathbf{0}$, and $\widehat{g}_{\mathbf{b}} = \mathbf{i}(\mathbf{b} + \mathbf{k}) \cdot \widehat{\mathbf{f}}_{\mathbf{b}} = -(\mathbf{b} + \mathbf{k}) \cdot ((\mathbf{b} + \mathbf{k}) \times \widehat{\mathbf{w}}_{\mathbf{b}}) = -\widehat{\mathbf{w}}_{\mathbf{b}} \cdot ((\mathbf{b} + \mathbf{k}) \times (\mathbf{b} + \mathbf{k})) = 0$ for all $\mathbf{b} \in \widehat{\Lambda}$. Hence, $\mathbf{v} = \mathbf{0}$ and $g = 0$. \square

Corollary 4.7. *For every vector $\mathbf{k} \in \mathbb{B}$ we have*

$$(\nabla + \mathbf{i}\mathbf{k}) \times \mathbf{H}_{\text{per}}(\text{curl}; \Omega) \perp (\nabla + \mathbf{i}\mathbf{k})H_{\text{per}}^1(\Omega),$$

where \perp denotes the orthogonality relation in $L^2(\Omega)^3$.

Proof. By Lemma 4.6 and Theorem 4.8 we have

$$\begin{aligned} \langle (\nabla + \mathbf{i}\mathbf{k}) \times \mathbf{w}, (\nabla + \mathbf{i}\mathbf{k})q \rangle_{\Omega} &= \int_{\Omega} \overline{(\nabla + \mathbf{i}\mathbf{k}) \times \mathbf{w}} \cdot (\nabla + \mathbf{i}\mathbf{k})q \\ &= \int_{\Omega} \overline{\mathbf{w}} \cdot (\nabla + \mathbf{i}\mathbf{k}) \times (\nabla + \mathbf{i}\mathbf{k})q = 0 \end{aligned}$$

for every $\mathbf{w} \in \mathbf{H}_{\text{per}}(\text{curl}; \Omega)$ and every $q \in H_{\text{per}}^1(\Omega)$. \square

Before we state our next theorem, we find it convenient to introduce the function space

$$L_{\diamond}^2(\Omega) := \left\{ f \in L^2(\Omega) \mid \int_{\Omega} f = 0 \right\}. \quad (4.24)$$

Clearly, $L_{\diamond}^2(\Omega)$ is a closed linear subspace of $L^2(\Omega)$ and thus a Hilbert space equipped with the inner product $\langle \cdot, \cdot \rangle_{\Omega}$. According to (4.5), we have

$$\widehat{f}_{\mathbf{0}} = \frac{1}{\text{meas}(\Omega)} \int_{\Omega} f \quad \text{for all } f \in L^2(\Omega).$$

Therefore, the function space $L_{\diamond}^2(\Omega)$ can be characterized as

$$L_{\diamond}^2(\Omega) = \{ f \in L^2(\Omega) \mid \widehat{f}_{\mathbf{0}} = 0 \}.$$

It is furthermore well-known that the Hilbert space $L_{\diamond}^2(\Omega)$ is isometrically linearly isomorphic to the linear factor space $L^2(\Omega)/\mathbb{C}$. From this follows the orthogonal space decomposition $L^2(\Omega) = L_{\diamond}^2(\Omega) \oplus \mathbb{C}$. The function spaces

$$\mathbf{H}_{\text{per}, \diamond}(\text{curl}; \Omega) := \mathbf{H}_{\text{per}}(\text{curl}; \Omega) \cap L_{\diamond}^2(\Omega)^3,$$

$$\begin{aligned}\mathbf{H}_{\text{per},\diamond}(\text{div};\Omega) &:= \mathbf{H}_{\text{per}}(\text{div};\Omega) \cap L_{\diamond}^2(\Omega)^3, \\ H_{\text{per},\diamond}^1(\Omega) &:= H_{\text{per}}^1(\Omega) \cap L_{\diamond}^2(\Omega),\end{aligned}$$

are closed linear subspaces of $\mathbf{H}_{\text{per}}(\text{curl};\Omega)$, $\mathbf{H}_{\text{per}}(\text{div};\Omega)$, and $H_{\text{per}}^1(\Omega)$, respectively.

Theorem 4.8. *The following assertions hold for every vector field $\mathbf{f} \in L^2(\Omega)^3$.*

- (a) *For every vector $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$ there exist uniquely determined functions $\boldsymbol{\psi} \in H_{\text{per}}^1(\Omega)^3$ and $p \in H_{\text{per}}^1(\Omega)$, as well as a constant $c > 0$, such that*

$$\begin{aligned}(\nabla + i\mathbf{k}) \times \boldsymbol{\psi} + (\nabla + i\mathbf{k})p &= \mathbf{f}, \\ (\nabla + i\mathbf{k}) \cdot \boldsymbol{\psi} &= 0, \\ \|\boldsymbol{\psi}\|_{1,\Omega} + \|p\|_{1,\Omega} &\leq c\|\mathbf{f}\|_{\Omega}.\end{aligned}$$

- (b) *There exist uniquely determined functions $\boldsymbol{\psi} \in H_{\text{per},\diamond}^1(\Omega)^3$ and $p \in H_{\text{per},\diamond}^1(\Omega)$, as well as a constant $c > 0$, such that*

$$\begin{aligned}\nabla \times \boldsymbol{\psi} + \nabla p + \widehat{\mathbf{f}}_{\mathbf{0}} &= \mathbf{f}, \\ \nabla \cdot \boldsymbol{\psi} &= 0, \\ \|\boldsymbol{\psi}\|_{1,\Omega} + \|p\|_{1,\Omega} + \sqrt{\text{meas}(\Omega)}|\widehat{\mathbf{f}}_{\mathbf{0}}| &\leq c\|\mathbf{f}\|_{\Omega}.\end{aligned}$$

Proof. A proof for (a) is given in Section 3 in [33], but can also be derived with the techniques that are presented in the following. In order to proof (b), we consider an arbitrary function $\mathbf{f} \in L^2(\Omega)^3$ with

$$\mathbf{f} = \sum_{\mathbf{b} \in \widehat{\Lambda}} \widehat{\mathbf{f}}_{\mathbf{b}} e^{i(\mathbf{b}, \cdot)}.$$

Notice that for every vector $\mathbf{b} \in \widehat{\Lambda} \setminus \{\mathbf{0}\}$ there exists a uniquely defined complex number $\widehat{q}_{\mathbf{b}} \in \mathbb{C}$ and a uniquely defined vector $\widehat{\boldsymbol{\varphi}}_{\mathbf{b}} \in (\text{span}\{\mathbf{b}\})^{\perp}$, such that $\widehat{\mathbf{f}}_{\mathbf{b}} = \mathbf{b}\widehat{q}_{\mathbf{b}} + \widehat{\boldsymbol{\varphi}}_{\mathbf{b}}$.

In Section 2.2 we introduced the cross product matrix $[\mathbf{z}]_{\times} \in \mathbb{C}^{3 \times 3}$ of a complex vector $\mathbf{z} \in \mathbb{C}^3$. For non-vanishing \mathbf{z} we showed that the image space of $[\mathbf{z}]_{\times}$ coincides with $(\text{span}\{\mathbf{z}\})^{\perp}$. It follows, that the pseudo-inverse $[\mathbf{z}]_{\times}^{\text{P}}$ of $[\mathbf{z}]_{\times}$ constitutes a linear isomorphism from $(\text{span}\{\mathbf{z}\})^{\perp}$ onto itself.

With the above definitions, we set $\widehat{p}_{\mathbf{b}} := -i\widehat{q}_{\mathbf{b}}$ and $\widehat{\boldsymbol{\psi}}_{\mathbf{b}} := -i[\mathbf{b}]_{\times}^{\text{P}}\widehat{\boldsymbol{\varphi}}_{\mathbf{b}}$. With this we obtain $\widehat{\mathbf{f}}_{\mathbf{b}} = i\mathbf{b}\widehat{p}_{\mathbf{b}} + i\mathbf{b} \times \widehat{\boldsymbol{\psi}}_{\mathbf{b}}$ for all $\mathbf{b} \in \widehat{\Lambda} \setminus \{\mathbf{0}\}$. Furthermore, since

$\text{im}([\mathbf{b}]_{\times}^{\text{P}}) = \text{im}([\mathbf{b}]_{\times}^{\text{T}}) = \text{im}([\mathbf{b}]_{\times})$, we have $\mathbf{b} \cdot \widehat{\boldsymbol{\psi}}_{\mathbf{b}} = \langle \mathbf{b}, \widehat{\boldsymbol{\psi}}_{\mathbf{b}} \rangle = \langle \mathbf{b}, \widehat{\boldsymbol{\varphi}}_{\mathbf{b}} \rangle = 0$ for all $\mathbf{b} \in \widehat{\Lambda}$. Next, we define the functions $p : \Omega \rightarrow \mathbb{C}$ and $\boldsymbol{\psi} : \Omega \rightarrow \mathbb{C}^3$ by

$$p := \sum_{\mathbf{b} \in \widehat{\Lambda} \setminus \{\mathbf{0}\}} \widehat{p}_{\mathbf{b}} e^{i\langle \mathbf{b}, \cdot \rangle}, \quad \boldsymbol{\psi} := \sum_{\mathbf{b} \in \widehat{\Lambda} \setminus \{\mathbf{0}\}} \widehat{\boldsymbol{\psi}}_{\mathbf{b}} e^{i\langle \mathbf{b}, \cdot \rangle}.$$

By Parseval's identity we have

$$\begin{aligned} \|p\|_{1,\Omega}^2 + \|\boldsymbol{\psi}\|_{1,\Omega}^2 &= \text{meas}(\Omega) \sum_{\mathbf{b} \in \widehat{\Lambda} \setminus \{\mathbf{0}\}} (1 + |\mathbf{b}|^2) (|\widehat{p}_{\mathbf{b}}|^2 + |\widehat{\boldsymbol{\psi}}_{\mathbf{b}}|^2) \\ &\leq 2 \text{meas}(\Omega) \sum_{\mathbf{b} \in \widehat{\Lambda} \setminus \{\mathbf{0}\}} (|\widehat{p}_{\mathbf{b}}|^2 |\mathbf{b}|^2 + |\mathbf{b}|^2 |\widehat{\boldsymbol{\psi}}_{\mathbf{b}}|^2) \\ &= 2 \text{meas}(\Omega) \sum_{\mathbf{b} \in \widehat{\Lambda} \setminus \{\mathbf{0}\}} (|\widehat{p}_{\mathbf{b}} \mathbf{b}|^2 + |\mathbf{b} \times \widehat{\boldsymbol{\psi}}_{\mathbf{b}}|^2) \\ &= 2 \text{meas}(\Omega) \sum_{\mathbf{b} \in \widehat{\Lambda} \setminus \{\mathbf{0}\}} (|\widehat{p}_{\mathbf{b}} \mathbf{b}|^2 + |\widehat{\boldsymbol{\varphi}}_{\mathbf{b}}|^2) \\ &= 2 \text{meas}(\Omega) \sum_{\mathbf{b} \in \widehat{\Lambda} \setminus \{\mathbf{0}\}} |\widehat{\mathbf{f}}_{\mathbf{b}}|^2 \\ &= 2 \|\mathbf{f}\|_{\Omega}^2 - 2 \text{meas}(\Omega) |\widehat{\mathbf{f}}_{\mathbf{0}}|^2. \end{aligned}$$

Since \mathbf{f} belongs to $L^2(\Omega)^3$, it follows that p and $\boldsymbol{\psi}$ belong to $H_{\text{per}}^1(\Omega)$ and $H_{\text{per}}^1(\Omega)^3$, respectively. According to (4.8)–(4.10) we also have that $\nabla \times \boldsymbol{\psi} + \nabla p + \widehat{\mathbf{f}}_{\mathbf{0}} = \mathbf{f}$, as well as $\nabla \cdot \boldsymbol{\psi} = 0$. Finally, the above inequality implies that $\|\boldsymbol{\psi}\|_{1,\Omega} + \|p\|_{1,\Omega} + \sqrt{\text{meas}(\Omega)} |\widehat{\mathbf{f}}_{\mathbf{0}}| \leq \sqrt{2} \|\mathbf{f}\|_{\Omega}$. \square

The functions $\boldsymbol{\psi}$ and p in Theorem 4.8 are commonly referred to as *vector potentials* and *scalar potentials* of the function \mathbf{f} . An additive decomposition of a vector field into a curl field of a vector potential and a gradient field of a scalar potential is called a *Helmholtz decomposition*. In analogy to this, Part (a) in Theorem 4.8 provides a Helmholtz decomposition with modified vector and scalar potentials $(\nabla + i\mathbf{k}) \times \boldsymbol{\psi}$ and $(\nabla + i\mathbf{k})p$, where $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$. It is also well-known, that for every vector field $\mathbf{f} \in L^2(\Omega)^3$ there exists a Helmholtz decomposition $\mathbf{f} = \nabla \times \boldsymbol{\psi} + \nabla p$, such that $\nabla \cdot \boldsymbol{\psi} = 0$, with uniquely determined functions $\boldsymbol{\psi} \in H_{\diamond}^1(\Omega)^3$ and $p \in H_{\diamond}^1(\Omega)$ (see e.g. Theorem and Remarks 3.3 in [38]). Part (b) of Theorem 4.8, however, states that such a decomposition with Λ -periodic vector and scalar potentials $\boldsymbol{\psi} \in H_{\text{per},\diamond}^1(\Omega)^3$ and $p \in H_{\text{per},\diamond}^1(\Omega)$ exists, if and only if $\mathbf{f} \in L_{\diamond}^2(\Omega)^3$.

Notice that Part (b) of Theorem 4.8 implies a fortiori that for every $\mathbf{f} \in L^2(\Omega)$ there exist functions $\boldsymbol{\psi} \in H_{\text{per}}^1(\Omega)^3$ and $p \in H_{\text{per}}^1(\Omega)$ with $\nabla \cdot \boldsymbol{\psi} = 0$, as well as a uniquely defined complex vector $\mathbf{z} \in \mathbb{C}^3$, such that $\mathbf{f} = \nabla \times \boldsymbol{\psi} + \nabla p + \mathbf{z}$. The

functions $\boldsymbol{\psi}$ and p are only uniquely determined modulo additive constants in this case. From Theorem 4.5 and Theorem 4.8 we hence deduce the following corollary.

Corollary 4.9. *For every vector $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$ we have the following orthogonal space decompositions.*

- (a) $L^2(\Omega)^3 = (\nabla + i\mathbf{k}) \times H_{\text{per}}^1(\Omega)^3 \oplus (\nabla + i\mathbf{k})H_{\text{per}}^1(\Omega).$
- (b) $L^2(\Omega)^3 = \nabla \times H_{\text{per},\diamond}^1(\Omega)^3 \oplus \nabla H_{\text{per},\diamond}^1(\Omega) \oplus \mathbb{C}^3.$
- (c) $L_{\diamond}^2(\Omega)^3 = \nabla \times H_{\text{per},\diamond}^1(\Omega)^3 \oplus \nabla H_{\text{per},\diamond}^1(\Omega).$

Note that Part (b) of Corollary 4.9 in particular states that the curl of any vector field in $H_{\text{per}}^1(\Omega)^3$ can only be constant, if it vanishes identically, and that the same holds true for the gradient of any function in $H_{\text{per}}^1(\Omega)$. Another important corollary is obtained from Lemma 4.6 and Theorem 4.8.

Corollary 4.10. *The following identities hold for every vector $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$.*

- (a) $\ker((\nabla + i\mathbf{k}) \times |_{\mathbf{H}_{\text{per}}(\text{curl};\Omega)}) = (\nabla + i\mathbf{k})H_{\text{per}}^1(\Omega).$
- (b) $\ker((\nabla + i\mathbf{k}) \cdot |_{\mathbf{H}_{\text{per}}(\text{div};\Omega)}) = (\nabla + i\mathbf{k}) \times \mathbf{H}_{\text{per}}(\text{curl};\Omega).$
- (c) $\ker(\nabla \times |_{\mathbf{H}_{\text{per}}(\text{curl};\Omega)}) = \nabla H_{\text{per}}^1(\Omega) \oplus \mathbb{C}^3.$
- (d) $\ker(\nabla \cdot |_{\mathbf{H}_{\text{per}}(\text{div};\Omega)}) = \nabla \times \mathbf{H}_{\text{per}}(\text{curl};\Omega) \oplus \mathbb{C}^3.$
- (e) $\ker(\nabla \times |_{\mathbf{H}_{\text{per},\diamond}(\text{curl};\Omega)}) = \nabla H_{\text{per},\diamond}^1(\Omega).$
- (f) $\ker(\nabla \cdot |_{\mathbf{H}_{\text{per},\diamond}(\text{div};\Omega)}) = \nabla \times \mathbf{H}_{\text{per},\diamond}(\text{curl};\Omega).$

In the literature the identities in Part (a) and (b) of Corollary 4.10 are often summarized in stating that for every vector $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$ the sequence

$$H_{\text{per}}^1(\Omega) \xrightarrow{(\nabla+i\mathbf{k})} \mathbf{H}_{\text{per}}(\text{curl};\Omega) \xrightarrow{(\nabla+i\mathbf{k})\times} \mathbf{H}_{\text{per}}(\text{div};\Omega) \xrightarrow{(\nabla+i\mathbf{k})\cdot} L^2(\Omega) \quad (4.25)$$

is exact (see e.g. [14], [33]). Here, “exact” simply means that the image of every operator in the sequence coincides with the kernel of the next operator. From Part (e) and (f) of Corollary 4.10 we further deduce that the sequence

$$H_{\text{per},\diamond}^1(\Omega) \xrightarrow{\nabla} \mathbf{H}_{\text{per},\diamond}(\text{curl};\Omega) \xrightarrow{\nabla\times} \mathbf{H}_{\text{per},\diamond}(\text{div};\Omega) \xrightarrow{\nabla\cdot} L_{\diamond}^2(\Omega) \quad (4.26)$$

is also exact. Both sequences are analogues of the *de Rham sequence*

$$H^1(\Omega) \xrightarrow{\nabla} \mathbf{H}(\text{curl};\Omega) \xrightarrow{\nabla\times} \mathbf{H}(\text{div};\Omega) \xrightarrow{\nabla\cdot} L^2(\Omega), \quad (4.27)$$

which is also known to be exact (see e.g. Proposition 8 and Theorem 8 in [20]). We remark that the sequences in (4.27) can be interpreted as a variant of de Rham's exact cochain complex of differential forms on \mathbb{R}^3 (see e.g. Section 1 in [16]). When interpreted this way, the function spaces $H^1(\Omega)$, $\mathbf{H}(\text{curl}; \Omega)$, $\mathbf{H}(\text{div}; \Omega)$, and $L^2(\Omega)$ represent certain classes of differential forms of order 0, 1, 2, and 3, respectively. The gradient, curl and divergence operators represent exterior derivative operators acting on the respective differential forms. Clearly, the sequences in (4.25) and (4.27) can be interpreted in the same way. Thus, all three sequences reveal certain geometric relationships between the respective function spaces.

4.3 The Weak Formulation

In this section we establish a weak formulation of the family of eigenvalue problems stated in Problem 4.1. It turns out that the weak formulation of each individual eigenvalue problem bears some similarity to the weak formulation of standard Laplace-type or second-order divergence type eigenvalue problems. In particular, one is able to prove coercivity and, for non-vanishing vectors $\mathbf{k} \in \mathbb{B}$, even ellipticity of the corresponding sesquilinear forms.

In order to establish the weak formulation, we define the complex Hilbert spaces

$$\mathbf{W} := \mathbf{H}_{\text{per}}(\text{curl}; \Omega), \quad (4.28)$$

$$Q := H_{\text{per}}^1(\Omega), \quad (4.29)$$

$$\mathbf{Z} := L^2(\Omega)^3, \quad (4.30)$$

as well as the real Banach space

$$\mathcal{E} := L^\infty(\Omega, \mathbb{R}), \quad (4.31)$$

and the function set

$$\mathcal{D} := \left\{ \rho \in \mathcal{E} \mid \exists \delta > 0 \text{ such that } \text{ess inf}_\Omega(\rho) \geq \delta \right\}. \quad (4.32)$$

The set \mathcal{D} consists of all essentially bounded functions, which are positive almost everywhere on Ω . Clearly, the function set \mathcal{D} is an open subset of \mathcal{E} . Furthermore, every function in \mathcal{D} satisfies the required properties of the coefficient ρ in Problem 4.1. Therefore, we define for every function $\rho \in \mathcal{D}$ and every vector $\mathbf{k} \in \mathbb{B}$ the sesquilinear forms $a_{\mathbf{k}}(\rho) : \mathbf{W} \times \mathbf{W} \rightarrow \mathbb{C}$, $b_{\mathbf{k}} : \mathbf{W} \times Q \rightarrow \mathbb{C}$, and $m : \mathbf{Z} \times \mathbf{Z} \rightarrow \mathbb{C}$ by

$$a_{\mathbf{k}}(\rho)(\mathbf{w}, \mathbf{v}) := \int_{\Omega} \rho \overline{(\nabla + i\mathbf{k}) \times \mathbf{w}} \cdot (\nabla + i\mathbf{k}) \times \mathbf{v} \quad \text{for all } \mathbf{w}, \mathbf{v} \in \mathbf{W}, \quad (4.33)$$

$$b_{\mathbf{k}}(\mathbf{f}, q) := \int_{\Omega} \bar{\mathbf{f}} \cdot (\nabla + i\mathbf{k})q \quad \text{for all } \mathbf{f} \in \mathbf{Z}, q \in Q, \quad (4.34)$$

$$m(\mathbf{f}, \mathbf{g}) := \int_{\Omega} \bar{\mathbf{f}} \cdot \mathbf{g} \quad \text{for all } \mathbf{f}, \mathbf{g} \in \mathbf{Z}. \quad (4.35)$$

Furthermore, we define for every vector $\mathbf{k} \in \mathbb{B}$ the linear function space

$$\mathbf{V}_{\mathbf{k}} := \{\mathbf{w} \in \mathbf{W} \mid b_{\mathbf{k}}(\mathbf{w}, q) = 0 \text{ for all } q \in Q\}. \quad (4.36)$$

Clearly, for every vector $\mathbf{k} \in \mathbb{B}$, the function space $\mathbf{V}_{\mathbf{k}}$ is a closed subspace of \mathbf{W} and thus a complex Hilbert space when equipped with the inner product $\langle \cdot, \cdot \rangle_{\text{curl}, \Omega}$. Each such Hilbert space $\mathbf{V}_{\mathbf{k}}$ can be characterized as follows.

Proposition 4.11. *For every vector $\mathbf{k} \in \mathbb{B}$ we have the identity*

$$\mathbf{V}_{\mathbf{k}} = \begin{cases} \mathbf{W} \cap ((\nabla + i\mathbf{k}) \times H_{\text{per}}^1(\Omega)^3) & \text{if } \mathbf{k} \neq \mathbf{0}, \\ \mathbf{W} \cap (\nabla \times H_{\text{per}}^1(\Omega)^3 \oplus \mathbb{C}^3) & \text{if } \mathbf{k} = \mathbf{0}. \end{cases}$$

Proof. First, we consider the case $\mathbf{k} \neq \mathbf{0}$. Let $\mathbf{w} \in \mathbf{W}$, with $\mathbf{w} = (\nabla + i\mathbf{k}) \times \boldsymbol{\psi}$ for some $\boldsymbol{\psi} \in H_{\text{per}}^1(\Omega)^3$. By Corollary 4.7 we then have

$$b_{\mathbf{k}}(\mathbf{w}, q) = \langle (\nabla + i\mathbf{k}) \times \boldsymbol{\psi}, (\nabla + i\mathbf{k})q \rangle_{\Omega} = 0 \quad \text{for all } q \in Q,$$

and hence $\mathbf{w} \in \mathbf{V}_{\mathbf{k}}$. Now, we choose a function $\mathbf{u} \in \mathbf{V}_{\mathbf{k}}$ arbitrarily. Then, by Theorem 4.8, there exist functions $\boldsymbol{\psi} \in H_{\text{per}}^1(\Omega)^3$ and $p \in Q$, such that $\mathbf{u} = (\nabla + i\mathbf{k}) \times \boldsymbol{\psi} + (\nabla + i\mathbf{k})p$. Furthermore, we have

$$\begin{aligned} \|(\nabla + i\mathbf{k})p\|_{\Omega}^2 &= \langle (\nabla + i\mathbf{k})p, (\nabla + i\mathbf{k})p \rangle_{\Omega} \\ &= \langle (\nabla + i\mathbf{k})p, (\nabla + i\mathbf{k})p \rangle_{\Omega} + \langle (\nabla + i\mathbf{k}) \times \boldsymbol{\psi}, (\nabla + i\mathbf{k})p \rangle_{\Omega} \\ &= b_{\mathbf{k}}((\nabla + i\mathbf{k})p, p) + b_{\mathbf{k}}((\nabla + i\mathbf{k}) \times \boldsymbol{\psi}, p) \\ &= b_{\mathbf{k}}(\mathbf{u}, p) = 0, \end{aligned}$$

which implies $(\nabla + i\mathbf{k})p = \mathbf{0}$. This completes the proof for the case $\mathbf{k} \neq \mathbf{0}$. The proof for the case $\mathbf{k} = \mathbf{0}$ is established analogously. Let $\mathbf{w}_0 \in \mathbf{W}$, with $\mathbf{w}_0 = \nabla \times \boldsymbol{\psi}_0 + \mathbf{z}_0$ for some function $\boldsymbol{\psi}_0 \in H_{\text{per}}^1(\Omega)^3$ and some vector $\mathbf{z}_0 \in \mathbb{C}^3$. By Theorem 4.8 and Corollary 4.7 we then obtain

$$b_0(\mathbf{w}_0, q) = \langle \nabla \times \boldsymbol{\psi}_0, \nabla q \rangle_{\Omega} + \langle \nabla \cdot \mathbf{z}_0, q \rangle_{\Omega} = 0 \quad \text{for all } q \in Q,$$

which implies that $\mathbf{w}_0 \in \mathbf{V}_0$. Choosing $\mathbf{u}_0 \in \mathbf{V}_0$ arbitrarily, we have by Theorem 4.8 that there exist functions $\boldsymbol{\psi}_0 \in H_{\text{per}}^1(\Omega)^3$ and $p_0 \in Q$, as well as a vector $\mathbf{z}_0 \in \mathbb{C}^3$, such that $\mathbf{u}_0 = \nabla \times \boldsymbol{\psi}_0 + \nabla p_0 + \mathbf{z}_0$. Furthermore, we have

$$\|\nabla p_0\|_{\Omega}^2 = \langle \nabla p_0, \nabla p_0 \rangle_{\Omega}$$

$$\begin{aligned}
&= \langle \nabla p_0, \nabla p_0 \rangle_\Omega + \langle \nabla \times \boldsymbol{\psi}_0, \nabla p_0 \rangle_\Omega + \langle \mathbf{z}_0, \nabla p_0 \rangle \\
&= b_{\mathbf{k}}(\nabla p_0, p_0) + b_{\mathbf{k}}(\nabla \times \boldsymbol{\psi}_0, p_0) + b_{\mathbf{k}}(\mathbf{z}_0, p_0) \\
&= b_{\mathbf{k}}(\mathbf{u}_0, p_0) = 0,
\end{aligned}$$

which implies $\nabla p_0 = \mathbf{0}$. □

In view of the eigenvalue problem stated in Problem 4.1 we have by Theorem 4.5 that each function space $\mathbf{V}_{\mathbf{k}}$ consists of those functions $\mathbf{u} \in \mathbf{W}$, which satisfy the divergence constraint $(\nabla + i\mathbf{k}) \cdot \mathbf{u} = 0$ in the weak sense. Another relevant aspect of the function spaces $\mathbf{V}_{\mathbf{k}}$ is the existence of the following norm estimate for $\mathbf{k} \neq \mathbf{0}$.

Lemma 4.12. *For every vector $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$ we have the norm estimate*

$$\|(\nabla + i\mathbf{k}) \times \mathbf{u}\|_\Omega \geq |\mathbf{k}| \|\mathbf{u}\|_\Omega \quad \text{for all } \mathbf{u} \in \mathbf{V}_{\mathbf{k}}.$$

Proof. Let $\mathbf{u} \in \mathbf{V}_{\mathbf{k}}$ be chosen arbitrarily. Then, according to Theorem 4.8 there exists a uniquely determined function $\boldsymbol{\psi} \in H_{\text{per}}^1(\Omega)^3$, such that $\mathbf{u} = (\nabla + i\mathbf{k}) \times \boldsymbol{\psi}$ and $(\nabla + i\mathbf{k}) \cdot \boldsymbol{\psi} = 0$. Note that the latter identity implies $(\mathbf{b} + \mathbf{k}) \cdot \widehat{\boldsymbol{\psi}}_{\mathbf{b}} = 0$ for all $\mathbf{b} \in \widehat{\Lambda}$, where $\widehat{\boldsymbol{\psi}}_{\mathbf{b}}$ denotes the respective Fourier coefficient of $\boldsymbol{\psi}$ as defined in (4.5). Using (2.4) and Parseval's identity, we obtain

$$\begin{aligned}
\frac{1}{\text{meas}(\Omega)} \|\mathbf{u}\|_\Omega^2 &= \frac{1}{\text{meas}(\Omega)} \|(\nabla + i\mathbf{k}) \times \boldsymbol{\psi}\|_\Omega^2 \\
&= \sum_{\mathbf{b} \in \widehat{\Lambda}} |(\mathbf{b} + \mathbf{k}) \times \widehat{\boldsymbol{\psi}}_{\mathbf{b}}|^2 \\
&= \sum_{\mathbf{b} \in \widehat{\Lambda}} (|\mathbf{b} + \mathbf{k}|^2 |\widehat{\boldsymbol{\psi}}_{\mathbf{b}}|^2 - |(\mathbf{b} + \mathbf{k}) \cdot \widehat{\boldsymbol{\psi}}_{\mathbf{b}}|^2) \\
&= \sum_{\mathbf{b} \in \widehat{\Lambda}} |\mathbf{b} + \mathbf{k}|^2 |\widehat{\boldsymbol{\psi}}_{\mathbf{b}}|^2.
\end{aligned}$$

Note that $|\mathbf{b} + \mathbf{k}| \geq |\mathbf{k}|$ for all $\mathbf{b} \in \widehat{\Lambda}$. For $\mathbf{k} = \mathbf{0}$ this assertion is trivial. For $\mathbf{k} \neq \mathbf{0}$ it follows from the fact that \mathbb{B} is the Wigner–Seitz cell of the reciprocal lattice (see Section 3.4). As was discussed in Section 2.5 the Wigner–Seitz cell of a lattice is point symmetric with respect to the origin and consists of all points, which are closer to the origin than to any other lattice point. Using the vector identity stated in (2.1), we find

$$\begin{aligned}
\frac{1}{\text{meas}(\Omega)} \|(\nabla + i\mathbf{k}) \times \mathbf{u}\|_\Omega^2 &= \frac{1}{\text{meas}(\Omega)} \|(\nabla + i\mathbf{k}) \times (\nabla + i\mathbf{k}) \times \boldsymbol{\psi}\|_\Omega^2 \\
&= \sum_{\mathbf{b} \in \widehat{\Lambda}} |(\mathbf{b} + \mathbf{k}) \times ((\mathbf{b} + \mathbf{k}) \times \widehat{\boldsymbol{\psi}}_{\mathbf{b}})|^2
\end{aligned}$$

$$\begin{aligned}
&= \sum_{\mathbf{b} \in \widehat{\Lambda}} |((\mathbf{b} + \mathbf{k}) \cdot \widehat{\boldsymbol{\psi}}_{\mathbf{b}})(\mathbf{b} + \mathbf{k}) - |\mathbf{b} + \mathbf{k}|^2 \widehat{\boldsymbol{\psi}}_{\mathbf{b}}|^2 \\
&= \sum_{\mathbf{b} \in \widehat{\Lambda}} |\mathbf{b} + \mathbf{k}|^4 |\widehat{\boldsymbol{\psi}}_{\mathbf{b}}|^2 \\
&\geq |\mathbf{k}|^2 \sum_{\mathbf{b} \in \widehat{\Lambda}} |\mathbf{b} + \mathbf{k}|^2 |\widehat{\boldsymbol{\psi}}_{\mathbf{b}}|^2,
\end{aligned}$$

which completes the proof. \square

Lemma 4.12 provides a *Poincaré–Friedrichs*-type inequality for the restriction of the modified curl operator $(\nabla + i\mathbf{k}) \times$ to the Hilbert space $\mathbf{V}_{\mathbf{k}}$ for every vector $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$. Since the modified curl operators also satisfy a Gårding-type inequality (see Lemma 4.4), we can state the following result.

Lemma 4.13. *For every vector $\mathbf{k} \in \mathbb{B}$ we have the norm estimate*

$$\|(\nabla + i\mathbf{k}) \times \mathbf{u}\|_{\Omega} \geq \sqrt{\gamma_0(\mathbf{k})} \|\mathbf{u}\|_{\text{curl}, \Omega} \quad \text{for all } \mathbf{u} \in \mathbf{V}_{\mathbf{k}}.$$

With the definitions in Lemma 4.4, the function $\gamma_0 : \mathbb{B} \rightarrow \mathbb{R}_{>0}$ is given by

$$\gamma_0(\mathbf{k}) := \frac{\alpha_0(\mathbf{k})|\mathbf{k}|^2}{|\mathbf{k}|^2 + \kappa_0(\mathbf{k})} \quad \text{for all } \mathbf{k} \in \mathbb{B}. \quad (4.37)$$

Proof. For $\mathbf{k} = \mathbf{0}$ the assertion is trivial. We hence turn our attention to the case, where $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$. Choosing a function $\mathbf{u} \in \mathbf{V}_{\mathbf{k}}$ arbitrarily, we have

$$\begin{aligned}
\alpha_0(\mathbf{k}) \|\mathbf{u}\|_{\text{curl}, \Omega}^2 &\leq \|(\nabla + i\mathbf{k}) \times \mathbf{u}\|_{\Omega}^2 + \kappa_0(\mathbf{k}) \|\mathbf{u}\|_{\Omega}^2 \\
&= \|(\nabla + i\mathbf{k}) \times \mathbf{u}\|_{\Omega}^2 + \frac{\kappa_0(\mathbf{k})}{|\mathbf{k}|^2} |\mathbf{k}|^2 \|\mathbf{u}\|_{\Omega}^2 \\
&\leq \left(1 + \frac{\kappa_0(\mathbf{k})}{|\mathbf{k}|^2}\right) \|(\nabla + i\mathbf{k}) \times \mathbf{u}\|_{\Omega}^2
\end{aligned}$$

by Lemma 4.4 and Lemma 4.12. The norm estimate is obtained by dividing the above inequality by the parenthesized term and taking the square root. \square

Lemma 4.13 and Lemma 4.3 imply, that $\|(\nabla + i\mathbf{k}) \times (\cdot)\|_{\Omega}$ constitutes a norm on $\mathbf{V}_{\mathbf{k}}$ for all $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$, which is equivalent to the norm $\|\cdot\|_{\text{curl}, \Omega}$. This is an analogon to the well-known fact that $\|\nabla(\cdot)\|_{\Omega}$, the seminorm on $H^1(\Omega)$, constitutes a norm on $H_0^1(\Omega)$, which is equivalent to the norm $\|\cdot\|_{1, \Omega}$. Note, however, that this analogon only holds for non-vanishing \mathbf{k} . In fact, one easily verifies that $\|\nabla \times \cdot\|_{\Omega}$ only constitutes a seminorm on $\mathbf{V}_{\mathbf{0}}$.

The following propositions state important properties of the sesquilinear forms $a_{\mathbf{k}}(\rho)$, which were defined in (4.33).

Proposition 4.14. *For every coefficient $\rho \in \mathcal{D}$, and every vector $\mathbf{k} \in \mathbb{B}$ the sesquilinear form $a_{\mathbf{k}}(\rho)$ is conjugate-symmetric, bounded, \mathbf{W} - \mathbf{Z} -coercive, and positive semidefinite. More precisely, we have*

$$\begin{aligned} a_{\mathbf{k}}(\rho)(\mathbf{w}, \mathbf{v}) &= \overline{a_{\mathbf{k}}(\rho)(\mathbf{v}, \mathbf{w})} && \text{for all } \mathbf{w}, \mathbf{v} \in \mathbf{W}, \\ |a_{\mathbf{k}}(\rho)(\mathbf{w}, \mathbf{v})| &\leq \beta(\rho, \mathbf{k}) \|\mathbf{w}\|_{\text{curl}, \Omega} \|\mathbf{v}\|_{\text{curl}, \Omega} && \text{for all } \mathbf{w}, \mathbf{v} \in \mathbf{W}, \\ a_{\mathbf{k}}(\rho)(\mathbf{w}, \mathbf{w}) &\geq \alpha(\rho, \mathbf{k}) \|\mathbf{w}\|_{\text{curl}, \Omega}^2 - \kappa(\rho, \mathbf{k}) \|\mathbf{w}\|_{\Omega}^2 && \text{for all } \mathbf{w} \in \mathbf{W}, \\ a_{\mathbf{k}}(\rho)(\mathbf{w}, \mathbf{w}) &\geq 0 && \text{for all } \mathbf{w} \in \mathbf{W}. \end{aligned}$$

With the definitions in Lemma 4.3 and Lemma 4.4 the functions $\alpha : \mathcal{D} \times \mathbb{B} \rightarrow \mathbb{R}_{>0}$, $\beta : \mathcal{D} \times \mathbb{B} \rightarrow \mathbb{R}_{>0}$, and $\kappa : \mathcal{D} \times \mathbb{B} \rightarrow \mathbb{R}_{>0}$ are given by

$$\alpha(\rho, \mathbf{k}) := \text{ess inf}_{\Omega}(\rho) \alpha_0(\mathbf{k}) \quad \text{for all } \rho \in \mathcal{D}, \mathbf{k} \in \mathbb{B}, \quad (4.38)$$

$$\beta(\rho, \mathbf{k}) := \|\rho\|_{\Omega, \infty} \beta_0(\mathbf{k}) \quad \text{for all } \rho \in \mathcal{D}, \mathbf{k} \in \mathbb{B}, \quad (4.39)$$

$$\kappa(\rho, \mathbf{k}) := \text{ess inf}_{\Omega}(\rho) \kappa_0(\mathbf{k}) \quad \text{for all } \rho \in \mathcal{D}, \mathbf{k} \in \mathbb{B}. \quad (4.40)$$

Proof. The conjugate-symmetry directly follows from the definition of the sesquilinear form $a_{\mathbf{k}}(\rho)$. Furthermore, choosing $\rho \in \mathcal{D}$, $\mathbf{k} \in \mathbb{B}$, and $\mathbf{w}, \mathbf{v} \in \mathbf{W}$ arbitrarily, we have

$$a_{\mathbf{k}}(\mathbf{w}, \mathbf{v}) \leq \|\rho\|_{\Omega, \infty} \|(\nabla + i\mathbf{k}) \times \mathbf{w}\|_{\Omega} \|(\nabla + i\mathbf{k}) \times \mathbf{v}\|_{\Omega}.$$

Boundedness thus follows from Lemma 4.3. Moreover, we have

$$a_{\mathbf{k}}(\mathbf{w}, \mathbf{w}) \geq \text{ess inf}_{\Omega}(\rho) \|(\nabla + i\mathbf{k}) \times \mathbf{w}\|_{\Omega}^2.$$

Hence, \mathbf{W} - \mathbf{Z} -coercivity follows from Lemma 4.4. The above inequality also implies the positive semidefiniteness of the sesquilinear form $a_{\mathbf{k}}$. \square

Proposition 4.15. *For every coefficient $\rho \in \mathcal{D}$ and every vector $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$ the sesquilinear form $a_{\mathbf{k}}(\rho)$ is $\mathbf{V}_{\mathbf{k}}$ -elliptic. More precisely, we have*

$$a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{u}) \geq \gamma(\rho, \mathbf{k}) \|\mathbf{u}\|_{\text{curl}, \Omega}^2 \quad \text{for all } \mathbf{u} \in \mathbf{V}_{\mathbf{k}}.$$

With the definitions in Lemma 4.13 the function $\gamma : \mathcal{D} \times \mathbb{B} \rightarrow \mathbb{R}_{>0}$ is given by

$$\gamma(\rho, \mathbf{k}) := \text{ess inf}_{\Omega}(\rho) \gamma_0(\mathbf{k}) \quad \text{for all } \rho \in \mathcal{D}, \mathbf{k} \in \mathbb{B}. \quad (4.41)$$

Proof. Given a coefficient $\rho \in \mathcal{D}$, a vector $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$, and a function $\mathbf{u} \in \mathbf{V}_{\mathbf{k}}$, we have

$$a_{\mathbf{k}}(\mathbf{u}, \mathbf{u}) \geq \text{ess inf}_{\Omega}(\rho) \|(\nabla + i\mathbf{k}) \times \mathbf{u}\|_{\Omega}^2.$$

The assertion hence follows from Lemma 4.13. \square

In the above Propositions 4.14 and 4.15 the terms *coercive* and *elliptic* were used as defined by Wloka (cf. Definitions 17.3 and 17.4 in [78]). Note that Proposition 4.15 implies that, for every $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$, the corresponding sesquilinear form $a_{\mathbf{k}}$ is positive definite on $\mathbf{V}_{\mathbf{k}}$.

Having stated some important properties of the function spaces $\mathbf{V}_{\mathbf{k}}$ and the sesquilinear forms $a_{\mathbf{k}}(\rho)$, we are finally in a position to state the weak formulation of Problem 4.1.

Problem 4.16. *Given a coefficient $\rho \in \mathcal{D}$, and a vector $\mathbf{k} \in \mathbb{B}$, find eigenvalues $\lambda \in \mathbb{C}$ and corresponding eigenfunctions $\mathbf{u} \in \mathbf{V}_{\mathbf{k}} \setminus \{\mathbf{0}\}$, such that*

$$a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{v}) = \lambda m(\mathbf{u}, \mathbf{v}) \quad \text{for all } \mathbf{v} \in \mathbf{V}_{\mathbf{k}}. \quad (4.42)$$

Note that the function spaces $\mathbf{V}_{\mathbf{k}}$ are sufficient as test spaces for the weak eigenvalue problem stated in Problem 4.16. Suppose that (4.42) holds for some coefficient $\rho \in \mathcal{D}$, some vector $\mathbf{k} \in \mathbb{B}$, some vector field $\mathbf{u} \in \mathbf{V}_{\mathbf{k}}$, and some complex number $\lambda \in \mathbb{C}$. Choosing an arbitrary test function $\mathbf{w} \in \mathbf{W}$, we have by Theorem 4.8 and Proposition 4.11 that there exist a vector field $\mathbf{v} \in \mathbf{V}_{\mathbf{k}}$ and a function $q \in Q$, such that

$$\mathbf{w} = \mathbf{v} + (\nabla + i\mathbf{k})q.$$

By Lemma 4.6 we then have that

$$a_{\mathbf{k}}(\rho)(\mathbf{u}, (\nabla + i\mathbf{k})q) = 0.$$

and hence

$$a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{w}) = a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{v}).$$

By definition of the sesquilinear forms m and $b_{\mathbf{k}}$, we also have that

$$m(\mathbf{u}, (\nabla + i\mathbf{k})q) = b_{\mathbf{k}}(\mathbf{u}, q) = 0,$$

since $\mathbf{u} \in \mathbf{V}_{\mathbf{k}}$. The latter identity implies that

$$m(\mathbf{u}, \mathbf{w}) = m(\mathbf{u}, \mathbf{v}).$$

It thus follows that (4.42) holds for some coefficient $\rho \in \mathcal{D}$, some vector $\mathbf{k} \in \mathbb{B}$, some vector field $\mathbf{u} \in \mathbf{V}_{\mathbf{k}}$, and some complex number $\lambda \in \mathbb{C}$, if and only if

$$a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{w}) = \lambda m(\mathbf{u}, \mathbf{w}) \quad \text{for all } \mathbf{w} \in \mathbf{W}.$$

It is therefore sufficient to choose $\mathbf{V}_{\mathbf{k}}$ as the test space in Problem 4.16.

4.4 Riesz–Schauder Theory

In this section we review the spectral theory for the weakly formulated family of eigenvalue problems stated in Problem 4.16. The theory we present here is standard and can be found in many textbooks (see e.g. Chapter 17 in [78]). We still discuss this theory here in detail, since some aspects are non-trivial for our particular family of eigenvalue problems.

In essence, we aim at applying the Spectral Theorem of Riesz–Schauder for compact, self-adjoint operators on infinite-dimensional Hilbert spaces. Our first step, therefore, is to devise such an operator. To this end, we define for every coefficient $\rho \in D$, and every vector $\mathbf{k} \in \mathbb{B}$ the conjugate-linear operators $A_{\mathbf{k}}(\rho) : \mathbf{V}_{\mathbf{k}} \rightarrow \mathbf{V}_{\mathbf{k}}^*$, $M_{\mathbf{k}} : \mathbf{Z} \rightarrow \mathbf{V}_{\mathbf{k}}^*$, and $I_{\mathbf{k}} : \mathbf{V}_{\mathbf{k}} \rightarrow \mathbf{Z}$ by

$$A_{\mathbf{k}}(\rho)\mathbf{u}[\mathbf{v}] := a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{v}) \quad \text{for all } \mathbf{u}, \mathbf{v} \in \mathbf{V}_{\mathbf{k}}, \quad (4.43)$$

$$M_{\mathbf{k}}\mathbf{f}[\mathbf{v}] := m(\mathbf{f}, \mathbf{v}) \quad \text{for all } \mathbf{f} \in \mathbf{Z}, \mathbf{v} \in \mathbf{V}_{\mathbf{k}}, \quad (4.44)$$

$$I_{\mathbf{k}}\mathbf{v} := \mathbf{v} \quad \text{for all } \mathbf{v} \in \mathbf{V}_{\mathbf{k}}. \quad (4.45)$$

Clearly, $I_{\mathbf{k}}$ is the identical embedding from $\mathbf{V}_{\mathbf{k}}$ into \mathbf{Z} for every $\mathbf{k} \in \mathbb{B}$. Note that the operator $M_{\mathbf{k}}$ is continuous for each $\mathbf{k} \in \mathbb{B}$ with

$$\|M_{\mathbf{k}}\|_{\mathbf{Z}, \mathbf{V}_{\mathbf{k}}^*} \leq 1. \quad (4.46)$$

Moreover, one easily verifies that

$$\ker(M_{\mathbf{k}}) = (\nabla + i\mathbf{k})Q \quad \text{for all } \mathbf{k} \in \mathbb{B}. \quad (4.47)$$

This implies, in particular, that the operators $M_{\mathbf{k}}$ are not injective. Restricting these operators to $\mathbf{V}_{\mathbf{k}}$, however, yields injective operators for all $\mathbf{k} \in \mathbb{B}$, since $(\nabla + i\mathbf{k})Q$ is \mathbf{Z} -orthogonal to $\mathbf{V}_{\mathbf{k}}$ according to Corollary 4.9 and Proposition 4.11.

With the above definitions, the weakly formulated eigenvalue problem stated in Problem 4.16 is equivalent to

Problem 4.17. *Given a coefficient $\rho \in D$, and a vector $\mathbf{k} \in \mathbb{B}$, find eigenvalues $\lambda \in \mathbb{C}$ and corresponding eigenfunctions $\mathbf{u} \in \mathbf{V}_{\mathbf{k}} \setminus \{\mathbf{0}\}$, such that*

$$A_{\mathbf{k}}(\rho)\mathbf{u} = \lambda M_{\mathbf{k}}\mathbf{u}.$$

Next, we define for every coefficient $\rho \in D$, and every vector $\mathbf{k} \in \mathbb{B}$ the conjugate-linear operator $L_{\mathbf{k}}(\rho) : \mathbf{V}_{\mathbf{k}} \rightarrow \mathbf{V}_{\mathbf{k}}^*$ by

$$L_{\mathbf{k}}(\rho) := A_{\mathbf{k}}(\rho) + \kappa(\rho, \mathbf{k}) M_{\mathbf{k}}|_{\mathbf{V}_{\mathbf{k}}}, \quad (4.48)$$

where $\kappa(\rho, \mathbf{k}) \in \mathbb{R}_{>0}$ is the constant defined in Proposition 4.14. It follows that every operator $L_{\mathbf{k}}(\rho)$ induces a conjugate-symmetric, bounded, and $\mathbf{V}_{\mathbf{k}}$ -elliptic sesquilinear form $l_{\mathbf{k}}(\rho) : \mathbf{V}_{\mathbf{k}} \times \mathbf{V}_{\mathbf{k}} \rightarrow \mathbb{C}$ by

$$l_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{v}) := L_{\mathbf{k}}(\rho)\mathbf{u}[\mathbf{v}] \quad \text{for all } \mathbf{u}, \mathbf{v} \in \mathbf{V}_{\mathbf{k}}. \quad (4.49)$$

More precisely, each sesquilinear form $l_{\mathbf{k}}(\rho)$ satisfies

$$\begin{aligned} l_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{v}) &= \overline{l_{\mathbf{k}}(\rho)(\mathbf{v}, \mathbf{u})} && \text{for all } \mathbf{u}, \mathbf{v} \in \mathbf{V}_{\mathbf{k}}, \\ |l_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{v})| &\leq \beta(\rho, \mathbf{k}) \|\mathbf{u}\|_{\text{curl}, \Omega} \|\mathbf{v}\|_{\text{curl}, \Omega} && \text{for all } \mathbf{u}, \mathbf{v} \in \mathbf{V}_{\mathbf{k}}, \\ l_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{u}) &\geq \alpha(\rho, \mathbf{k}) \|\mathbf{u}\|_{\text{curl}, \Omega}^2 && \text{for all } \mathbf{u} \in \mathbf{V}_{\mathbf{k}}, \end{aligned}$$

where $\alpha(\rho, \mathbf{k})$ and $\beta(\rho, \mathbf{k})$ were defined in Proposition 4.14. The Lemma of Lax–Milgram then implies that every operator $L_{\mathbf{k}}(\rho)$ has a conjugate-linear, continuous inverse $L_{\mathbf{k}}(\rho)^{-1} : \mathbf{V}_{\mathbf{k}}^* \rightarrow \mathbf{V}_{\mathbf{k}}$ with

$$\|L_{\mathbf{k}}(\rho)^{-1}\|_{\mathbf{V}_{\mathbf{k}}^*, \mathbf{V}_{\mathbf{k}}} \leq \frac{1}{\alpha(\rho, \mathbf{k})} \quad \text{for all } \rho \in \mathcal{D}, \mathbf{k} \in \mathbb{B}. \quad (4.50)$$

Therefore, we are able to define for every coefficient $\rho \in \mathcal{D}$, and every vector $\mathbf{k} \in \mathbb{B}$ the linear operator $G_{\mathbf{k}}(\rho) : \mathbf{Z} \rightarrow \mathbf{Z}$ by

$$G_{\mathbf{k}}(\rho) := I_{\mathbf{k}} L_{\mathbf{k}}(\rho)^{-1} M_{\mathbf{k}} \quad \text{for all } \mathbf{k} \in \mathbb{B}, \quad (4.51)$$

In the literature, the operator $G_{\mathbf{k}}(\rho)$ is sometimes referred to as the *Green solution operator* for $L_{\mathbf{k}}(\rho)$ (see e.g. Definition 17.5 in [78]). The following theorem states that the identical embeddings $I_{\mathbf{k}}$ from $\mathbf{V}_{\mathbf{k}}$ into \mathbf{Z} are compact operators.

Theorem 4.18. *For every vector $\mathbf{k} \in \mathbb{B}$ the complex Hilbert space $\mathbf{V}_{\mathbf{k}}$ is compactly embedded in \mathbf{Z} .*

Proof. The following proof is an adaptation of the proof of Theorem 2.1 in [74]. Let $\{\mathbf{u}^{(n)}\}_{n \in \mathbb{N}}$ be a sequence in $\mathbf{V}_{\mathbf{k}}$, which is bounded in $\mathbf{V}_{\mathbf{k}}$ by a constant $c_1 > 0$. By Theorem 4.8 and Proposition 4.11 there exists a sequence $\{\boldsymbol{\psi}^{(n)}\}_{n \in \mathbb{N}}$ in $H_{\text{per}}^1(\Omega)^3$, a sequence $\{\mathbf{z}_n\}_{n \in \mathbb{N}}$ in \mathbb{C}^3 , and a constant $c > 0$, such that

$$\begin{aligned} (\nabla + i\mathbf{k}) \times \boldsymbol{\psi}^{(n)} + \mathbf{z}_n &= \mathbf{u}^{(n)} && \text{for all } n \in \mathbb{N}, \\ \|\boldsymbol{\psi}^{(n)}\|_{1, \Omega} + \sqrt{\text{meas}(\Omega)} |\mathbf{z}_n| &\leq c \|\mathbf{u}^{(n)}\|_{\Omega} && \text{for all } n \in \mathbb{N}. \end{aligned}$$

Since $\|\mathbf{u}^{(n)}\|_{\Omega} \leq \|\mathbf{u}^{(n)}\|_{\text{curl}, \Omega} \leq c_1$ for all $n \in \mathbb{N}$, the sequences $\{\boldsymbol{\psi}^{(n)}\}_{n \in \mathbb{N}}$ and $\{\mathbf{z}_n\}_{n \in \mathbb{N}}$ are bounded in $H^1(\Omega)^3$ and \mathbb{C}^3 , respectively. According to the Theorem of Rellich–Kondrachov (see e.g. Theorem 6.3 in [1]) the identical embedding from $H^1(\Omega)^3$ into \mathbf{Z} is compact. Therefore, there exists a subsequence $\{\boldsymbol{\psi}^{(1, n)}\}_{n \in \mathbb{N}}$

of $\{\boldsymbol{\psi}^{(n)}\}_{n \in \mathbb{N}}$, which converges in \mathbf{Z} . According to the Theorem of Heine–Borel there also exists a subsequence $\{\mathbf{z}_{(2,n)}\}_{n \in \mathbb{N}}$ of $\{\mathbf{z}_{(1,n)}\}_{n \in \mathbb{N}}$, which converges in \mathbb{C}^3 . We now consider the corresponding subsequence $\{\mathbf{u}^{(2,n)}\}_{n \in \mathbb{N}}$ of $\{\mathbf{u}^{(n)}\}_{n \in \mathbb{N}}$. To simplify notations, we define $\mathbf{u}^{[m,n]} := \mathbf{u}^{(2,m)} - \mathbf{u}^{(2,n)}$, $\boldsymbol{\psi}^{[m,n]} := \boldsymbol{\psi}^{(2,m)} - \boldsymbol{\psi}^{(2,n)}$, and $\mathbf{z}^{[m,n]} := \mathbf{z}_{(2,m)} - \mathbf{z}_{(2,n)}$ for all $m, n \in \mathbb{N}$. By Proposition 4.3 and Theorem 4.5, we obtain

$$\begin{aligned}
& \|\mathbf{u}^{(2,m)} - \mathbf{u}^{(2,n)}\|_{\Omega}^2 \\
&= \int_{\Omega} \overline{\mathbf{u}^{[m,n]}} \cdot \mathbf{u}^{[m,n]} \\
&= \int_{\Omega} \overline{\mathbf{u}^{[m,n]}} \cdot ((\nabla + i\mathbf{k}) \times \boldsymbol{\psi}^{[m,n]} + \mathbf{z}^{[m,n]}) \\
&= \int_{\Omega} \overline{(\nabla + i\mathbf{k}) \times \mathbf{u}^{[m,n]}} \cdot \boldsymbol{\psi}^{[m,n]} + \int_{\Omega} \overline{\mathbf{u}^{[m,n]}} \cdot \mathbf{z}^{[m,n]} \\
&\leq \|(\nabla + i\mathbf{k}) \times \mathbf{u}^{[m,n]}\|_{\Omega} \|\boldsymbol{\psi}^{[m,n]}\|_{\Omega} + \sqrt{\text{meas}(\Omega)} \|\mathbf{u}^{[m,n]}\|_{\Omega} |\mathbf{z}^{[m,n]}| \\
&\leq \|\mathbf{u}^{[m,n]}\|_{\text{curl}, \Omega} (\beta_0(\mathbf{k}) \|\boldsymbol{\psi}^{[m,n]}\|_{\Omega} + \sqrt{\text{meas}(\Omega)} |\mathbf{z}^{[m,n]}|) \\
&\leq 2c_1 (\beta_0(\mathbf{k}) \|\boldsymbol{\psi}^{[m,n]}\|_{\Omega} + \sqrt{\text{meas}(\Omega)} |\mathbf{z}^{[m,n]}|) \\
&\leq c_2 (\|\boldsymbol{\psi}^{(2,m)} - \boldsymbol{\psi}^{(2,n)}\|_{\Omega} + |\mathbf{z}_{(2,m)} - \mathbf{z}_{(2,n)}|) \rightarrow 0 \quad \text{as } m, n \rightarrow \infty,
\end{aligned}$$

where $c_2 := 2c_1 \max\{\beta_0(\mathbf{k}), \sqrt{\text{meas}(\Omega)}\}$. Hence, $\{\mathbf{u}^{(2,n)}\}_{n \in \mathbb{N}}$ is a Cauchy sequence in \mathbf{Z} and therefore convergent in \mathbf{Z} . \square

The following proposition establishes important properties of the operators

Proposition 4.19. *The following assertions hold for every coefficient $\rho \in D$ and every vector $\mathbf{k} \in \mathbb{B}$.*

- (a) *The operator $G_{\mathbf{k}}(\rho)$ is compact, self-adjoint and positive semidefinite.*
- (b) *$\ker(G_{\mathbf{k}}(\rho)) = \mathbf{V}_{\mathbf{k}}^{\perp}$, where $\mathbf{V}_{\mathbf{k}}^{\perp}$ denotes the \mathbf{Z} -orthogonal complement of $\mathbf{V}_{\mathbf{k}}$.*

Proof. In the following we use the fact that $\langle \mathbf{f}, \mathbf{v} \rangle_{\Omega} = M_{\mathbf{k}} \mathbf{f}[\mathbf{v}]$ for all $\mathbf{f} \in \mathbf{Z}$ and all $\mathbf{v} \in \mathbf{V}_{\mathbf{k}}$. Furthermore, we use that fact that $I_{\mathbf{k}}$ is the identity on $\mathbf{V}_{\mathbf{k}}$.

(a) Compactness directly follows from the definition of $G_{\mathbf{k}}(\rho)$. In order to prove self-adjointness and positive semi-definiteness, we choose two functions $\mathbf{f}, \mathbf{g} \in \mathbf{Z}$ arbitrarily. Then, there exist uniquely determined functions $\mathbf{u}, \mathbf{w} \in \mathbf{V}_{\mathbf{k}}$, such that $L_{\mathbf{k}}(\rho)\mathbf{u} = M_{\mathbf{k}}\mathbf{f}$ and $L_{\mathbf{k}}(\rho)\mathbf{w} = M_{\mathbf{k}}\mathbf{g}$. We hence obtain

$$\langle G_{\mathbf{k}}(\rho)\mathbf{f}, \mathbf{g} \rangle_{\Omega} = \langle I_{\mathbf{k}} L_{\mathbf{k}}(\rho)^{-1} M_{\mathbf{k}}\mathbf{f}, \mathbf{g} \rangle_{\Omega} = \langle L_{\mathbf{k}}(\rho)^{-1} M_{\mathbf{k}}\mathbf{f}, \mathbf{g} \rangle_{\Omega} = \langle \mathbf{u}, \mathbf{g} \rangle_{\Omega}$$

$$\begin{aligned}
&= \overline{\langle \mathbf{g}, \mathbf{u} \rangle_\Omega} = \overline{M_{\mathbf{k}} \mathbf{g}[\mathbf{u}]} = \overline{L_{\mathbf{k}}(\rho) \mathbf{w}[\mathbf{u}]} = L_{\mathbf{k}}(\rho) \mathbf{u}[\mathbf{w}] = M_{\mathbf{k}} \mathbf{f}[\mathbf{w}] \\
&= \langle \mathbf{f}, \mathbf{w} \rangle_\Omega = \langle \mathbf{f}, L_{\mathbf{k}}(\rho)^{-1} M_{\mathbf{k}} \mathbf{g} \rangle_\Omega = \langle \mathbf{f}, I_{\mathbf{k}} L_{\mathbf{k}}(\rho)^{-1} M_{\mathbf{k}} \mathbf{g} \rangle_\Omega \\
&= \langle \mathbf{f}, G_{\mathbf{k}}(\rho) \mathbf{g} \rangle_\Omega,
\end{aligned}$$

which proves the self-adjointedness of $G_{\mathbf{k}}(\rho)$. Similarly, we deduce the inequality

$$\langle G_{\mathbf{k}}(\rho) \mathbf{f}, \mathbf{f} \rangle_\Omega = L_{\mathbf{k}}(\rho) \mathbf{u}[\mathbf{u}] \geq \alpha(\rho, \mathbf{k}) \|\mathbf{u}\|_{\text{curl}, \Omega} \geq 0,$$

which proves the positive semidefiniteness of $G_{\mathbf{k}}(\rho)$.

(b) For every $\mathbf{f} \in \mathbf{Z}$, and every $\mathbf{v} \in \mathbf{V}_{\mathbf{k}}$ we have the identity

$$\begin{aligned}
\langle \mathbf{f}, \mathbf{v} \rangle_\Omega &= M_{\mathbf{k}} \mathbf{f}[\mathbf{v}] = (L_{\mathbf{k}}(\rho) L_{\mathbf{k}}(\rho)^{-1} M_{\mathbf{k}} \mathbf{f})[\mathbf{v}] = L_{\mathbf{k}}(\rho) (L_{\mathbf{k}}(\rho)^{-1} M_{\mathbf{k}} \mathbf{f})[\mathbf{v}] \\
&= L_{\mathbf{k}}(\rho) (I_{\mathbf{k}} L_{\mathbf{k}}(\rho)^{-1} M_{\mathbf{k}} \mathbf{f})[\mathbf{v}] = L_{\mathbf{k}}(\rho) (G_{\mathbf{k}}(\rho) \mathbf{f})[\mathbf{v}].
\end{aligned}$$

Since $L_{\mathbf{k}}(\rho)$ is a bijective operator, $\langle \mathbf{f}, \mathbf{v} \rangle_\Omega = 0$ holds for all $\mathbf{v} \in \mathbf{V}_{\mathbf{k}}$, if and only if $G_{\mathbf{k}}(\rho) \mathbf{f} = 0$. Hence, $\ker(G_{\mathbf{k}}(\rho)) = \mathbf{V}_{\mathbf{k}}^\perp$. \square

Since the operators $G_{\mathbf{k}}(\rho)$ are compact, they are subject to the Spectral Theorem of Riesz–Schauder (see e.g. Section 12.1 in [78]). Hence, the following statements hold for every coefficient $\rho \in \mathcal{D}$, and every vector $\mathbf{k} \in \mathbb{B}$. First, we have the orthogonal decomposition

$$\mathbf{Z} = \ker(G_{\mathbf{k}}(\rho)) \oplus \overline{\text{im}(G_{\mathbf{k}}(\rho))}.$$

By Part (b) of Proposition 4.19 we hence have

$$\mathbf{V}_{\mathbf{k}} = \overline{\text{im}(G_{\mathbf{k}}(\rho))}. \quad (4.52)$$

Furthermore, since the operator $G_{\mathbf{k}}(\rho)$ is self-adjointed and positive semidefinite, we have that there exists a sequence $\{\mu_j(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$, which consists of positive, real eigenvalues of $G_{\mathbf{k}}(\rho)$, and which converges to zero. In the following we shall assume that the eigenvalues of $G_{\mathbf{k}}(\rho)$ are enumerated, such that

$$\mu_1(\rho, \mathbf{k}) \geq \mu_2(\rho, \mathbf{k}) \geq \mu_3(\rho, \mathbf{k}) \geq \dots$$

The sequence of corresponding eigenfunctions $\{\mathbf{u}_j(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$ forms a complete \mathbf{Z} -orthogonal system of $\mathbf{V}_{\mathbf{k}}$. Here, as in the following, we denote by $\mathbf{u}_j(\rho, \mathbf{k})$ the eigenfunction that corresponds to the j -th largest eigenvalue of $G_{\mathbf{k}}(\rho)$, i.e.,

$$G_{\mathbf{k}}(\rho) \mathbf{u}_j(\rho, \mathbf{k}) = \mu_j(\rho, \mathbf{k}) \mathbf{u}_j(\rho, \mathbf{k}). \quad (4.53)$$

The eigenspace of each eigenvalue is finite-dimensional. The spectrum of the operator $G_{\mathbf{k}}(\rho)$ consists precisely of the elements of the sequence $\{\mu_j(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$ and

zero, which is also an eigenvalue of $G_{\mathbf{k}}(\rho)$. Finally, we have that the operator $G_{\mathbf{k}}(\rho)$ acts on \mathbf{Z} as

$$G_{\mathbf{k}}(\rho)\mathbf{f} := \sum_{j=1}^{\infty} \mu_j(\rho, \mathbf{k}) \frac{\langle \mathbf{u}_j(\rho, \mathbf{k}), \mathbf{f} \rangle_{\Omega}}{\|\mathbf{u}_j(\rho, \mathbf{k})\|_{\Omega}^2} \mathbf{u}_j(\rho, \mathbf{k}) \quad \text{for all } \mathbf{f} \in \mathbf{Z}.$$

The j -th largest eigenvalue can be characterized according to the *Min-Max Principle* as

$$\mu_j(\rho, \mathbf{k}) = \max_{\substack{\mathbf{U} \sqsubseteq \mathbf{Z} \\ \dim \mathbf{U} = j}} \min_{\mathbf{f} \in \mathbf{U} \setminus \{\mathbf{0}\}} \frac{\langle G_{\mathbf{k}}(\rho)\mathbf{f}, \mathbf{f} \rangle_{\Omega}}{\|\mathbf{f}\|_{\Omega}^2} \quad \text{for all } j \in \mathbb{N}, \quad (4.54)$$

where \sqsubseteq denotes the linear subspace relation (see Theorem XIII.1 in [64]). For convenience, we define for every coefficient $\rho \in \mathcal{D}$, every vector $\mathbf{k} \in \mathbb{B}$, and every index $j \in \mathbb{N}$ the linear function space

$$\mathbf{V}_{\mathbf{k},j}(\rho) := \text{span}\{\mathbf{u}_1(\rho, \mathbf{k}), \dots, \mathbf{u}_j(\rho, \mathbf{k})\}. \quad (4.55)$$

Then, one can show that the maxima and minima in (4.54) are attained precisely for $\mathbf{U} = \mathbf{V}_{\mathbf{k},j}(\rho)$. An equivalent characterization of the j -th largest eigenvalue is given by *Rayleigh's principle*, which reads

$$\mu_j(\rho, \mathbf{k}) = \max_{\substack{\mathbf{f} \in \mathbf{Z} \setminus \{\mathbf{0}\} \\ P_{\mathbf{k},j-1}(\rho)\mathbf{f} = \mathbf{0}}} \frac{\langle G_{\mathbf{k}}(\rho)\mathbf{f}, \mathbf{f} \rangle_{\Omega}}{\|\mathbf{f}\|_{\Omega}^2}, \quad (4.56)$$

Here, $P_{\mathbf{k},j}(\rho)$ is the linear, \mathbf{Z} -orthogonal projection operator from \mathbf{Z} onto the linear function space $\mathbf{V}_{\mathbf{k},j}(\rho)$. for every coefficient $\rho \in \mathcal{D}$, every vector $\mathbf{k} \in \mathbb{B}$, and every index $j \in \mathbb{N}$ the operator $P_{\mathbf{k},j}(\rho) : \mathbf{Z} \rightarrow \mathbf{V}_{\mathbf{k}}(\rho)$ is given by

$$P_{\mathbf{k},j}(\rho)\mathbf{f} := \sum_{i=1}^j \frac{\langle \mathbf{u}_i(\rho, \mathbf{k}), \mathbf{f} \rangle_{\Omega}}{\|\mathbf{u}_i(\rho, \mathbf{k})\|_{\Omega}^2} \mathbf{u}_i(\rho, \mathbf{k}) \quad \text{for all } \mathbf{f} \in \mathbf{Z}. \quad (4.57)$$

Now, let $\mu_j(\rho, \mathbf{k}) > 0$ be a positive eigenvalue of $G_{\mathbf{k}}(\rho)$, and let $\mathbf{u}_j(\rho, \mathbf{k})$ be a corresponding eigenfunction, i.e.,

$$G_{\mathbf{k}}(\rho)\mathbf{u}_j(\rho, \mathbf{k}) = \mu_j(\rho, \mathbf{k})\mathbf{u}_j(\rho, \mathbf{k}). \quad (4.58)$$

According to (4.48) and (4.51) we can rewrite (4.58) equivalently as

$$\left(\frac{1}{\mu_j(\rho, \mathbf{k})} - \kappa(\rho, \mathbf{k}) \right) M_{\mathbf{k}}\mathbf{u}_j(\rho, \mathbf{k}) = A_{\mathbf{k}}(\rho)\mathbf{u}_j(\rho, \mathbf{k}). \quad (4.59)$$

Since the sequence of eigenfunctions $\{\mathbf{u}_j(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$ of the operator $G_{\mathbf{k}}(\rho)$ form a complete orthogonal system of $\mathbf{V}_{\mathbf{k}}$, it follows that each eigenpair of $G_{\mathbf{k}}(\rho)$ corresponds to a solution of the eigenvalue problem stated in Problem 4.17 and thus to a solution of the weak eigenvalue problem stated in Problem 4.16. In particular, the elements of the sequence $\{\lambda_j(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$, defined by

$$\lambda_j(\rho, \mathbf{k}) := \frac{1}{\mu_j(\rho, \mathbf{k})} - \kappa(\rho, \mathbf{k}) \quad \text{for all } j \in \mathbb{N}, \quad (4.60)$$

are the eigenvalues. In summary, we have the following proposition.

Proposition 4.20. *For every coefficient $\rho \in \mathcal{D}$ and every vector $\mathbf{k} \in \mathbb{B}$ there exists an increasing sequence $\{\lambda_j(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$ of non-negative, real numbers tending to ∞ , which consists of the eigenvalues of the weakly formulated eigenvalue problem (4.42). The sequence of corresponding eigenfunctions $\{\mathbf{u}_j(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$ forms a complete, \mathbf{Z} -orthogonal system of $\mathbf{V}_{\mathbf{k}}$. The eigenspace of each eigenvalue is finite-dimensional. The j -th smallest eigenvalue can be characterized as*

$$\lambda_j(\rho, \mathbf{k}) = \min_{\substack{\mathbf{U} \subseteq \mathbf{V}_{\mathbf{k}} \\ \dim \mathbf{U} = j}} \max_{\mathbf{u} \in \mathbf{U} \setminus \{0\}} \frac{a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{u})}{m(\mathbf{u}, \mathbf{u})} \quad \text{for all } j \in \mathbb{N}. \quad (4.61)$$

or equivalently as

$$\lambda_j(\rho, \mathbf{k}) = \min_{\substack{\mathbf{u} \in \mathbf{V}_{\mathbf{k}} \\ P_{\mathbf{k}, j-1}(\rho)\mathbf{u} = 0}} \frac{a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{u})}{m(\mathbf{u}, \mathbf{u})} \quad \text{for all } j \in \mathbb{N}. \quad (4.62)$$

Zero is an eigenvalue, if and only if $\mathbf{k} = 0$.

Proof. The existence of the sequences $\{\lambda_j(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$ and $\{\mathbf{u}_j(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$ follows from the equivalence of (4.58) and (4.59). Notice also, that every eigenvalue $\lambda_j(\rho, \mathbf{k})$ is an eigenvalue of the weakly formulated eigenvalue problem (4.42).

In order to show (4.61), we define for every $\rho \in \mathcal{D}$ and every $\mathbf{k} \in \mathbb{B}$ the operators $F_{\mathbf{k}}(\rho) : \mathbf{V}_{\mathbf{k}} \rightarrow \mathbf{V}_{\mathbf{k}}$ by $F_{\mathbf{k}} := G_{\mathbf{k}}|_{\mathbf{V}_{\mathbf{k}}}$, as well as the operator $F_{\mathbf{k}}^{1/2}(\rho) : \mathbf{V}_{\mathbf{k}} \rightarrow \mathbf{V}_{\mathbf{k}}$ by

$$F_{\mathbf{k}}^{1/2}(\rho) \mathbf{v} := \sum_{j=1}^{\infty} \sqrt{\mu_j(\rho, \mathbf{k})} \frac{\langle \mathbf{u}_j(\rho, \mathbf{k}), \mathbf{v} \rangle_{\Omega}}{\|\mathbf{u}_j(\rho, \mathbf{k})\|_{\Omega}^2} \mathbf{u}_j(\rho, \mathbf{k}) \quad \text{for all } \mathbf{u} \in \mathbf{V}_{\mathbf{k}}.$$

One easily verifies that the continuous, linear operators $F_{\mathbf{k}}(\rho)$ and $F_{\mathbf{k}}^{1/2}(\rho)$ are symmetric with respect to the inner product $\langle \cdot, \cdot \rangle_{\Omega}$ and continuously invertible for all $\rho \in \mathcal{D}$ and all $\mathbf{k} \in \mathbb{B}$. The operators also satisfy $F_{\mathbf{k}}(\rho) = F_{\mathbf{k}}^{1/2}(\rho)F_{\mathbf{k}}^{1/2}(\rho)$ and $G_{\mathbf{k}}(\rho)F_{\mathbf{k}}(\rho)^{-1} = \text{id}_{\mathbf{V}_{\mathbf{k}}}$, where $\text{id}_{\mathbf{V}_{\mathbf{k}}}$ denotes the identity on $\mathbf{V}_{\mathbf{k}}$, for all $\rho \in \mathcal{D}$ and

all $\mathbf{k} \in \mathbb{B}$. Choosing an arbitrary function $\mathbf{u} \in \mathbf{V}_{\mathbf{k}}$, we define $\mathbf{v} := F_{\mathbf{k}}(\rho)^{-1}\mathbf{u}$ and $\mathbf{w} := L_{\mathbf{k}}(\rho)^{-1}M_{\mathbf{k}}\mathbf{v}$ for some $\rho \in \mathcal{D}$ and some $\mathbf{k} \in \mathbb{B}$. We then find that

$$\mathbf{w} = L_{\mathbf{k}}(\rho)^{-1}M_{\mathbf{k}}F_{\mathbf{k}}(\rho)^{-1}\mathbf{u} = I_{\mathbf{k}}L_{\mathbf{k}}(\rho)^{-1}M_{\mathbf{k}}F_{\mathbf{k}}(\rho)^{-1}\mathbf{u} = G_{\mathbf{k}}(\rho)F_{\mathbf{k}}(\rho)^{-1}\mathbf{u} = \mathbf{u},$$

and hence,

$$\begin{aligned} \langle F_{\mathbf{k}}(\rho)^{-1}\mathbf{u}, \mathbf{u} \rangle_{\Omega} &= \langle \mathbf{v}, F_{\mathbf{k}}(\rho)\mathbf{v} \rangle_{\Omega} = \langle \mathbf{v}, G_{\mathbf{k}}(\rho)\mathbf{v} \rangle_{\Omega} = \langle \mathbf{v}, I_{\mathbf{k}}L_{\mathbf{k}}(\rho)^{-1}M_{\mathbf{k}}\mathbf{v} \rangle_{\Omega} \\ &= \langle \mathbf{v}, L_{\mathbf{k}}(\rho)^{-1}M_{\mathbf{k}}\mathbf{v} \rangle_{\Omega} = \langle \mathbf{v}, \mathbf{w} \rangle_{\Omega} = M_{\mathbf{k}}\mathbf{v}[\mathbf{w}] = L_{\mathbf{k}}(\rho)\mathbf{w}[\mathbf{w}] \\ &= L_{\mathbf{k}}(\rho)\mathbf{u}[\mathbf{u}]. \end{aligned}$$

Notice that this identity can be established for all functions $\mathbf{u} \in \mathbf{V}_{\mathbf{k}}$. According to the Min-Max Principle (4.54), we have

$$\begin{aligned} \frac{1}{\mu_j(\rho, \mathbf{k})} &= \left(\max_{\substack{U \subseteq \mathbf{Z} \\ \dim U = j}} \min_{\mathbf{f} \in U \setminus \{0\}} \frac{\langle G_{\mathbf{k}}(\rho)\mathbf{f}, \mathbf{f} \rangle_{\Omega}}{\|\mathbf{f}\|_{\Omega}^2} \right)^{-1} \\ &= \min_{\substack{U \subseteq \mathbf{Z} \\ \dim U = j}} \max_{\mathbf{f} \in U \setminus \{0\}} \frac{\|\mathbf{f}\|_{\Omega}^2}{\langle G_{\mathbf{k}}(\rho)\mathbf{f}, \mathbf{f} \rangle_{\Omega}} \\ &= \min_{\substack{U \subseteq \mathbf{V}_{\mathbf{k}} \\ \dim U = j}} \max_{\mathbf{u} \in U \setminus \{0\}} \frac{\|\mathbf{u}\|_{\Omega}^2}{\langle F_{\mathbf{k}}(\rho)\mathbf{u}, \mathbf{u} \rangle_{\Omega}} \\ &= \min_{\substack{U \subseteq \mathbf{V}_{\mathbf{k}} \\ \dim U = j}} \max_{\mathbf{u} \in U \setminus \{0\}} \frac{\langle \mathbf{u}, \mathbf{u} \rangle_{\Omega}}{\langle F_{\mathbf{k}}^{1/2}(\rho)\mathbf{u}, F_{\mathbf{k}}^{1/2}(\rho)\mathbf{u} \rangle_{\Omega}} \\ &= \min_{\substack{U \subseteq \mathbf{V}_{\mathbf{k}} \\ \dim U = j}} \max_{\mathbf{u} \in U \setminus \{0\}} \frac{\langle F_{\mathbf{k}}^{1/2}(\rho)^{-1}\mathbf{u}, F_{\mathbf{k}}^{1/2}(\rho)^{-1}\mathbf{u} \rangle_{\Omega}}{\langle \mathbf{u}, \mathbf{u} \rangle_{\Omega}} \\ &= \min_{\substack{U \subseteq \mathbf{V}_{\mathbf{k}} \\ \dim U = j}} \max_{\mathbf{u} \in U \setminus \{0\}} \frac{\langle F_{\mathbf{k}}(\rho)^{-1}\mathbf{u}, \mathbf{u} \rangle_{\Omega}}{\langle \mathbf{u}, \mathbf{u} \rangle_{\Omega}} \\ &= \min_{\substack{U \subseteq \mathbf{V}_{\mathbf{k}} \\ \dim U = j}} \max_{\mathbf{u} \in U \setminus \{0\}} \frac{L_{\mathbf{k}}(\rho)\mathbf{u}[\mathbf{u}]}{M_{\mathbf{k}}\mathbf{u}[\mathbf{u}]} \\ &= \min_{\substack{U \subseteq \mathbf{V}_{\mathbf{k}} \\ \dim U = j}} \max_{\mathbf{u} \in U \setminus \{0\}} \frac{a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{u})}{m(\mathbf{u}, \mathbf{u})} + \kappa(\rho, \mathbf{k}) \quad \text{for all } j \in \mathbb{N}. \end{aligned}$$

This establishes the characterization given in (4.61). The characterization given in (4.62) can be shown analogously using Rayleigh's principle (4.56).

Finally, we show that zero is an eigenvalue if and only if $\mathbf{k} \in \mathbb{B}$ vanishes. According to Proposition 4.15 the sesquilinear forms $a_{\mathbf{k}}(\rho)$ are $\mathbf{V}_{\mathbf{k}}$ -elliptic for all

$\rho \in \mathcal{D}$ and all $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$. Therefore, zero cannot be an eigenvalue, if $\mathbf{k} \neq \mathbf{0}$. By Proposition 4.11, however, the three-dimensional space of constant functions from Ω into \mathbb{C}^3 is a linear subspace of \mathbf{V}_0 . One easily verifies that a non-vanishing constant function $\mathbf{u} : \Omega \rightarrow \mathbb{C}^3$ satisfies $a_0(\mathbf{u}, \mathbf{u}) = 0$. Hence, zero is an eigenvalue for $\mathbf{k} = \mathbf{0}$. \square

Proposition 4.20 makes a conclusive statement about the existence of Bloch modes in photonic crystals. Recall that Problem 4.16 was motivated by the family of constrained eigenvalue problems (3.36), which arose from the Bloch ansatz for the magnetic field of a time-harmonic wave propagating inside a photonic crystal (see Section 4.1). The coefficient ρ represents the restriction to Ω of the crystal's reciprocal relative electric permittivity function $1/\varepsilon_r$. The vector $\mathbf{k} \in \mathbb{B}$ coincides with the quasimomentum vector of a Bloch mode. Proposition 4.20 now states that in every three-dimensional photonic crystal there exists a countable set of Bloch modes for each quasimomentum vector \mathbf{k} . The frequency of each Bloch mode is determined by the corresponding eigenvalue λ of (4.42) via $\omega = c_0\sqrt{\lambda}$. The Bloch modes for a given quasimomentum vector \mathbf{k} form a complete \mathbf{Z} -orthogonal system.

4.5 Auchmuty's Principle

The identities (4.61) and (4.62) in Proposition 4.20 allow us to characterize the eigenvalues of Problem 4.16 in terms of the sesquilinear forms $a_{\mathbf{k}}(\rho)$, as defined by (4.33). The first identity is a variant of the Min-Max Principle, the second one is a variant of Rayleigh's principle. In this section we introduce yet another, however less well-known, characterization principle which goes back to works of Auchmuty (cf. [8], [9]). We present this principle here, because it will play an important role in proving the existence of optimal solutions for photonic band gap optimization problems.

In essence, Auchmuty's principle stems from a duality principle for certain types of optimization problems. Before we can make this notion precise, we need to introduce some fundamental concepts of convex analysis. The concepts we present here are discussed in detail in standard textbooks on convex analysis such as [35] or [65]. In the following we assume that X and Y are real Banach spaces, and that X^* and Y^* denote their normed duals. A convex functional from X into the set $\overline{\mathbb{R}} := \mathbb{R} \cup \{-\infty, \infty\}$ is called *proper* if it is not identical to $-\infty$ or ∞ . Given a proper, lower semicontinuous, convex functional $f : X \rightarrow \overline{\mathbb{R}}$, the so-called *convex conjugate* functional $f^* : X^* \rightarrow \overline{\mathbb{R}}$ of f is defined by

$$f^*(p) := \sup_{x \in X} (p(x) - f(x)) \quad \text{for all } p \in X^*, \quad (4.63)$$

where X^* denotes the normed dual of X . One can show that f^* is also a proper, lower semicontinuous, convex functional.

Given two proper, lower semicontinuous, convex functionals $f : X \rightarrow \overline{\mathbb{R}}$ and $g : Y \rightarrow \overline{\mathbb{R}}$, we define the functional $F : X \rightarrow \overline{\mathbb{R}}$ by

$$F(x) := f(x) - g(Tx) \quad \text{for all } x \in X,$$

where $T : X \rightarrow Y$ is a bijective, continuous linear operator. Notice that the functional F is not convex in general. Our aim is to find the infimum

$$\inf_{x \in X} F(x). \quad (4.64)$$

One observes that F can be characterized as

$$F(x) = \inf_{q \in Y^*} \mathcal{L}(x, q) \quad \text{for all } x \in X, \quad (4.65)$$

where the functional $\mathcal{L} : X \times Y^* \rightarrow \overline{\mathbb{R}}$ is given by

$$\mathcal{L}(x, q) := f(x) + g^*(-q) + q(Tx) \quad \text{for all } x \in X, q \in Y^*.$$

The functional \mathcal{L} is called a *Lagrangian of type II* associated with F . Here, “of type II” means that the functional F is determined by the infima of \mathcal{L} with respect to the Lagrange parameter q according to (4.65). In his work, Auchmuty shows that the following *duality principle* holds (see Theorem 3.3 in [8]),

$$\inf_{x \in X} \inf_{q \in Y^*} \mathcal{L}(x, q) = \inf_{q \in Y^*} \inf_{x \in X} \mathcal{L}(x, q). \quad (4.66)$$

One can show, that the functional $G : Y^* \rightarrow \overline{\mathbb{R}}$, given by

$$G(q) := \inf_{x \in X} \mathcal{L}(x, q) \quad \text{for all } q \in Y^*,$$

is characterized as

$$G(q) = g^*(-q) - f^*(-T^*q) \quad \text{for all } q \in Y^*.$$

Here, T^* denotes the *dual operator* of T in the sense of functional analysis. We hence find, that the infimum (4.64) coincides with

$$\inf_{q \in Y^*} G(q). \quad (4.67)$$

Moreover, the duality principle (4.66) can be rewritten as

$$\inf_{x \in X} (f(x) - g(Tx)) = \inf_{q \in Y^*} (g^*(-q) - f^*(-T^*q)). \quad (4.68)$$

We remark that similar duality principles are well-known for convex functionals $F : X \rightarrow \overline{\mathbb{R}}$, which are of the form $F = f + g \circ T$ (see e.g. Chapter 3, Section 4 in [35]).

An interesting application of the duality principle stated in (4.68) is the following characterization principle for the eigenvalues of the weakly formulated eigenvalue problem (4.42).

Proposition 4.21. *Under the assumptions of Proposition 4.20, we have*

$$\frac{-1}{2(\lambda_j(\rho, \mathbf{k}) + \kappa(\rho, \mathbf{k}))} = \min_{\mathbf{u} \in \mathbf{V}_{\mathbf{k}}} \left(\frac{1}{2} (a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{u}) + \kappa(\rho, \mathbf{k})m(\mathbf{u}, \mathbf{u})) - \|\mathbf{u} - P_{\mathbf{k}, j-1}(\rho)\mathbf{u}\|_{\Omega} \right) \quad \text{for all } j \in \mathbb{N}.$$

Proof. This proof is an adaptation of the line of arguments laid out in [9]. In essence, we aim at taking advantage of the duality principle stated in (4.68). To this end, we define the real Hilbert space $\mathcal{Z} := L^2(\Omega, \mathbb{R})^3$. Furthermore, we define the isometric, linear isomorphism $\Phi : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathcal{Z}$ by

$$\Phi(\mathbf{f}_1, \mathbf{f}_2) := \mathbf{f}_1 + i\mathbf{f}_2 \quad \text{for all } \mathbf{f}_1, \mathbf{f}_2 \in \mathcal{Z}.$$

Choosing an arbitrary coefficient $\rho \in \mathcal{D}$ and a vector $\mathbf{k} \in \mathbb{B}$, we define for every index $j \in \mathbb{N}$ the set

$$\mathcal{K}_j := \{(\mathbf{f}_1, \mathbf{f}_2) \in \mathcal{Z} \times \mathcal{Z} \mid P_{\mathbf{k}, j-1}(\rho)\Phi(\mathbf{f}_1, \mathbf{f}_2) = \mathbf{0}, \|\Phi(\mathbf{f}_1, \mathbf{f}_2)\|_{\Omega} \leq 1\}.$$

One observes that for every $j \in \mathbb{N}$ the set \mathcal{K}_j is a closed, convex subset of $\mathcal{Z} \times \mathcal{Z}$. Next, we define for every index $j \in \mathbb{N}$ the functional $\psi_j : \mathcal{Z} \times \mathcal{Z} \rightarrow \overline{\mathbb{R}}$ as well as the functional $\sigma : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$ by

$$\psi_j(\mathbf{f}_1, \mathbf{f}_2) := \begin{cases} 0 & \text{if } (\mathbf{f}_1, \mathbf{f}_2) \in \mathcal{K}_j, \\ \infty & \text{else} \end{cases} \quad \text{for all } \mathbf{f}_1, \mathbf{f}_2 \in \mathcal{Z},$$

$$\sigma(\mathbf{f}_1, \mathbf{f}_2) := \frac{1}{2} \langle G_{\mathbf{k}}(\rho)\Phi(\mathbf{f}_1, \mathbf{f}_2), \Phi(\mathbf{f}_1, \mathbf{f}_2) \rangle_{\Omega} \quad \text{for all } \mathbf{f}_1, \mathbf{f}_2 \in \mathcal{Z}.$$

Clearly, the functionals ψ_j and σ are proper, lower semicontinuous, and convex on $\mathcal{Z} \times \mathcal{Z}$. In particular, the functional ψ_j is the so-called *indicator function* of the set \mathcal{K}_j in the sense of convex analysis (see e.g. page 28 in [65]).

Next, we determine the convex conjugate functionals of ψ_j for $j \in \mathbb{N}$. Notice that the Rieszian isomorphism $R : \mathcal{Z} \times \mathcal{Z} \rightarrow (\mathcal{Z} \times \mathcal{Z})^*$ is given by

$$R(\mathbf{f}_1, \mathbf{f}_2)[\mathbf{g}_1, \mathbf{g}_2] := \langle \mathbf{f}_1, \mathbf{g}_1 \rangle_{\Omega} + \langle \mathbf{f}_2, \mathbf{g}_2 \rangle_{\Omega} \quad \text{for all } \mathbf{f}_1, \mathbf{f}_2, \mathbf{g}_1, \mathbf{g}_2 \in \mathcal{Z}.$$

Choosing two functions $\mathbf{f}_1, \mathbf{f}_2 \in \mathcal{Z}$ arbitrarily, we set $\mathbf{f} := \Phi(\mathbf{f}_1, \mathbf{f}_2)$. Given an arbitrary index $j \in \mathbb{N}$, we then obtain

$$\begin{aligned}
(\psi_j^* \circ R)(\mathbf{f}_1, \mathbf{f}_2) &= \sup_{(\mathbf{g}_1, \mathbf{g}_2) \in \mathcal{Z} \times \mathcal{Z}} \left(\langle \mathbf{f}_1, \mathbf{g}_1 \rangle_\Omega + \langle \mathbf{f}_2, \mathbf{g}_2 \rangle_\Omega - \psi_j(\mathbf{g}_1, \mathbf{g}_2) \right) \\
&= \sup_{\mathbf{g} \in \mathcal{Z}} \left(\operatorname{Re}(\langle \mathbf{f}, \mathbf{g} \rangle_\Omega) - \psi_j(\operatorname{Re}(\mathbf{g}), \operatorname{Im}(\mathbf{g})) \right) \\
&= \sup \left\{ \operatorname{Re}(\langle \mathbf{f}, \mathbf{g} \rangle_\Omega) \mid \mathbf{g} \in \mathcal{Z}, \|\mathbf{g}\|_\Omega \leq 1, P_{\mathbf{k}, j-1}(\rho)\mathbf{g} = \mathbf{0} \right\} \\
&= \sup \left\{ \operatorname{Re}(\langle \mathbf{f}, \mathbf{g} \rangle_\Omega) \mid \mathbf{g} \in \mathcal{Z}, \|\mathbf{g}\|_\Omega = 1, P_{\mathbf{k}, j-1}(\rho)\mathbf{g} = \mathbf{0} \right\} \\
&= \sup_{\mathbf{h} \in \mathcal{Z}} \frac{\operatorname{Re}(\langle \mathbf{f}, \mathbf{h} - P_{\mathbf{k}, j-1}(\rho)\mathbf{h} \rangle_\Omega)}{\|\mathbf{h} - P_{\mathbf{k}, j-1}(\rho)\mathbf{h}\|_\Omega} \\
&= \sup_{\mathbf{h} \in \mathcal{Z}} \frac{\operatorname{Re}(\langle \mathbf{f} - P_{\mathbf{k}, j-1}(\rho)\mathbf{f}, \mathbf{h} - P_{\mathbf{k}, j-1}(\rho)\mathbf{h} \rangle_\Omega)}{\|\mathbf{h} - P_{\mathbf{k}, j-1}(\rho)\mathbf{h}\|_\Omega}.
\end{aligned}$$

According to the Cauchy–Schwarz inequality, the maximum is attained for $\mathbf{h} = \mathbf{f}$. Hence, for every index $j \in \mathbb{N}$ the convex conjugate functional of ψ_j is given by

$$(\psi_j^* \circ R)(\mathbf{f}_1, \mathbf{f}_2) = \|\Phi(\mathbf{f}_1, \mathbf{f}_2) - P_{\mathbf{k}, j-1}(\rho)\Phi(\mathbf{f}_1, \mathbf{f}_2)\|_\Omega \quad \text{for all } \mathbf{f}_1, \mathbf{f}_2 \in \mathcal{Z}.$$

We remark that ψ_j^* coincides with the so-called *support function* of \mathcal{K}_j in the sense of convex analysis (see e.g. page 28 in [65]), i.e.,

$$\psi_j^*(q) = \sup_{(\mathbf{f}_1, \mathbf{f}_2) \in \mathcal{K}_j} q(\mathbf{f}_1, \mathbf{f}_2) \quad \text{for all } q \in (\mathcal{Z} \times \mathcal{Z})^*.$$

Now, we determine the convex conjugate functional of σ . Again, we choose two functions $\mathbf{f}_1, \mathbf{f}_2 \in \mathcal{Z}$ arbitrarily, and set $\mathbf{f} := \Phi(\mathbf{f}_1, \mathbf{f}_2)$. If $\mathbf{f} \in \ker(G_{\mathbf{k}}(\rho))$ and $\mathbf{f} \neq \mathbf{0}$, we find that

$$\begin{aligned}
(\sigma^* \circ R)(\mathbf{f}_1, \mathbf{f}_2) &= \sup_{(\mathbf{g}_1, \mathbf{g}_2) \in \mathcal{Z} \times \mathcal{Z}} \left(\langle \mathbf{f}_1, \mathbf{g}_1 \rangle_\Omega + \langle \mathbf{f}_2, \mathbf{g}_2 \rangle_\Omega - \sigma(\mathbf{g}_1, \mathbf{g}_2) \right) \\
&= \sup_{\mathbf{g} \in \mathcal{Z}} \left(\operatorname{Re}(\langle \mathbf{f}, \mathbf{g} \rangle_\Omega) - \frac{1}{2} \langle G_{\mathbf{k}}(\rho)\mathbf{g}, \mathbf{g} \rangle_\Omega \right) \\
&\geq \operatorname{Re}(\langle \mathbf{f}, s\mathbf{f} \rangle_\Omega) - \frac{1}{2} \langle sG_{\mathbf{k}}(\rho)\mathbf{f}, s\mathbf{f} \rangle_\Omega \\
&= s\|\mathbf{f}\|_\Omega^2,
\end{aligned}$$

where $s > 0$ is an arbitrary positive number. Since s can be chosen arbitrarily, we conclude that $(\sigma^* \circ R)(\mathbf{f}_1, \mathbf{f}_2) = \infty$ if $\mathbf{f} \in \ker(G_{\mathbf{k}}(\rho))$. Next, we assume that

$\mathbf{f} \in \text{im}(G_{\mathbf{k}}(\rho))$. Then, there exists a function $\mathbf{g}_{\mathbf{f}} \in \mathbf{Z}$, such that $G_{\mathbf{k}}(\rho)\mathbf{g}_{\mathbf{f}} = \mathbf{f}$. Moreover, we have that the functional $\varphi_{\mathbf{f}} : \mathbf{Z} \rightarrow \mathbb{R}$, defined by

$$\varphi_{\mathbf{f}}(\mathbf{g}) := \text{Re}(\langle \mathbf{f}, \mathbf{g} \rangle_{\Omega}) - \frac{1}{2} \langle G_{\mathbf{k}}(\rho)\mathbf{g}, \mathbf{g} \rangle \quad \text{for all } \mathbf{g} \in \mathbf{Z},$$

is maximized on the affine subspace $\mathbf{g}_{\mathbf{f}} + \ker(G_{\mathbf{k}}(\rho))$. To see this, one can simply compute the first and second Fréchet derivatives $D\varphi_{\mathbf{f}}$ and $D^2\varphi_{\mathbf{f}}$ of the functional $\varphi_{\mathbf{f}}$. Then, it is easy to verify that $D\varphi_{\mathbf{f}}(\mathbf{g}) = 0$, if and only if $\mathbf{g} \in \mathbf{g}_{\mathbf{f}} + \ker(G_{\mathbf{k}}(\rho))$, and that $D^2\varphi_{\mathbf{f}}(\mathbf{v}, \mathbf{v}) < 0$ for all $\mathbf{v} \in \ker(G_{\mathbf{k}}(\rho))^{\perp}$. Furthermore, since

$$\mathbf{f} = G_{\mathbf{k}}(\rho)\mathbf{g}_{\mathbf{f}} = I_{\mathbf{k}}L_{\mathbf{k}}(\rho)^{-1}M_{\mathbf{k}}\mathbf{g}_{\mathbf{f}} = L_{\mathbf{k}}(\rho)^{-1}M_{\mathbf{k}}\mathbf{g}_{\mathbf{f}},$$

we have that

$$\langle \mathbf{f}, \mathbf{g}_{\mathbf{f}} \rangle_{\Omega} = M_{\mathbf{k}}\mathbf{g}_{\mathbf{f}}[\mathbf{f}] = L_{\mathbf{k}}(\rho)\mathbf{f}[\mathbf{f}],$$

and hence

$$(\sigma^* \circ R)(\mathbf{f}_1, \mathbf{f}_2) = \sup_{\mathbf{g} \in \mathbf{Z}} \varphi_{\mathbf{f}}(\mathbf{g}) = \varphi_{\mathbf{f}}(\mathbf{g}_{\mathbf{f}}) = \frac{1}{2} \langle \mathbf{f}, \mathbf{g}_{\mathbf{f}} \rangle_{\Omega} = \frac{1}{2} L_{\mathbf{k}}(\rho)\mathbf{f}[\mathbf{f}].$$

Since the image space of $G_{\mathbf{k}}(\rho)$ is dense in $\mathbf{V}_{\mathbf{k}}$ according to (4.52), we conclude that the convex conjugate functional of σ is given by

$$(\sigma^* \circ R)(\mathbf{f}_1, \mathbf{f}_2) = \begin{cases} \frac{1}{2} L_{\mathbf{k}}(\rho)\Phi(\mathbf{f}_1, \mathbf{f}_2)[\Phi(\mathbf{f}_1, \mathbf{f}_2)] & \text{if } \Phi(\mathbf{f}_1, \mathbf{f}_2) \in \mathbf{V}_{\mathbf{k}}, \\ \infty & \text{else} \end{cases}$$

for all $\mathbf{f}_1, \mathbf{f}_2 \in \mathbf{Z}$. According to Rayleigh's principle (4.56), we have

$$-\frac{1}{2} \mu_j(\rho, \mathbf{k}) = \min_{(\mathbf{f}_1, \mathbf{f}_2) \in \mathbf{Z} \times \mathbf{Z}} \left(\psi_j(\mathbf{f}_1, \mathbf{f}_2) - \sigma(\mathbf{f}_1, \mathbf{f}_2) \right) \quad \text{for all } j \in \mathbb{N}.$$

The assertion thus follows from (4.60) and the duality principle stated in (4.68). \square

By reviewing the line of arguments in the previous Section 4.4, one realizes that the operator $\kappa(\rho, \mathbf{k})M_{\mathbf{k}}|_{\mathbf{V}_{\mathbf{k}}}$ was added to $A_{\mathbf{k}}(\rho)$ in order to obtain a continuously invertible operator $L_{\mathbf{k}}(\rho) = A_{\mathbf{k}}(\rho) + \kappa(\rho, \mathbf{k})M_{\mathbf{k}}|_{\mathbf{V}_{\mathbf{k}}}$ for all $\rho \in \mathcal{D}$ and all $\mathbf{k} \in \mathbb{B}$. The continuous invertibility of $L_{\mathbf{k}}(\rho)$ followed from the ellipticity of the associated sesquilinear forms $l_{\mathbf{k}}(\rho) = a_{\mathbf{k}}(\rho) + \kappa(\rho, \mathbf{k})m$. However, Proposition 4.15 states that the sesquilinear forms $a_{\mathbf{k}}(\rho)$ are also elliptic for non-vanishing \mathbf{k} . Hence, the spectral theory in Section 4.4 can also be established with $\kappa(\rho, \mathbf{k})$ being replaced by zero, provided that \mathbf{k} does not vanish. From this we deduce the following variant of Auchmuty's principle.

Corollary 4.22. *Under the assumptions of Proposition 4.20, the following identity holds for all $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$,*

$$\frac{-1}{2\lambda_j(\rho, \mathbf{k})} = \min_{\mathbf{u} \in \mathbf{V}_{\mathbf{k}}} \left(\frac{1}{2} a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{u}) - \|\mathbf{u} - P_{\mathbf{k}, j-1}(\rho)\mathbf{u}\|_{\Omega} \right) \quad \text{for all } j \in \mathbb{N}.$$

Note that for $j = 1$ the assertion of Corollary 4.22 can be obtained without using the duality principle (4.68). Given a coefficient $\rho \in \mathcal{D}$ and a vector $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$, one simply defines the functional $f : \mathbf{V}_{\mathbf{k}} \rightarrow \mathbb{R}$ by

$$f(\mathbf{u}) := \frac{1}{2} a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{u}) - \|\mathbf{u}\|_{\Omega} \quad \text{for all } \mathbf{u} \in \mathbf{V}_{\mathbf{k}}.$$

Clearly, the functional f is continuous on $\mathbf{V}_{\mathbf{k}}$. Since $f(\mathbf{0}) = 0$ and $f(\theta\mathbf{v}) \rightarrow \infty$ as $\theta \rightarrow \infty$ for all $\mathbf{v} \in \mathbf{V}_{\mathbf{k}}$, we deduce that f admits a global minimum on $\mathbf{V}_{\mathbf{k}}$. The Fréchet-derivative of f is given by

$$Df(\mathbf{u})[\mathbf{v}] = a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{v}) - \frac{m(\mathbf{u}, \mathbf{v})}{\|\mathbf{u}\|_{\Omega}} \quad \text{for all } \mathbf{u}, \mathbf{v} \in \mathbf{V}_{\mathbf{k}}.$$

From this we see that any stationary point $\mathbf{u}_* \in \mathbf{V}_{\mathbf{k}}$ of f solves the weak eigenvalue equation (4.42) for the eigenvalue $\lambda_* = 1/\|\mathbf{u}_*\|_{\Omega}$. Therefore, we have

$$f(\mathbf{u}_*) = -\frac{1}{2\lambda_*}.$$

It follows that the function values of f at its stationary points are given by $-1/(2\lambda)$, where λ is a eigenvalue of the weak eigenvalue equation (4.42). The minimal function value at a critical point is therefore given by $-1/(2\lambda_1(\rho, \mathbf{k}))$, i.e.,

$$\frac{-1}{2\lambda_1(\rho, \mathbf{k})} = \min_{\mathbf{u} \in \mathbf{V}_{\mathbf{k}}} \left(\frac{1}{2} a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{u}) - \|\mathbf{u}\|_{\Omega} \right).$$

We conclude this section with the remark that Auchmuty's principle, as stated in Proposition 4.21 or Corollary 4.22, can be established for every coercive or elliptic sesquilinear form on a complex Hilbert space, respectively.

4.6 Photonic Band Structures and Band Diagrams

In Section 4.4, we studied the existence of eigensolutions of the weakly formulated family of constrained eigenvalue problems given by Problem 4.16. We found that for every coefficient $\rho \in \mathcal{D}$ and every vector $\mathbf{k} \in \mathbb{B}$ there exists a countable number

of eigensolutions. Provided that the coefficient ρ represents the medium structure of a three-dimensional photonic crystal, each eigensolution corresponds to a Bloch mode of that crystal with quasimomentum vector \mathbf{k} . In this section we relate the eigensolutions of Problem 4.16 to the so-called band structures of photonic crystals.

According to Proposition 4.20, there exists an increasing sequence $\{\lambda_j(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$ of non-negative eigenvalues of (4.42) for every $\rho \in \mathcal{D}$ and every $\mathbf{k} \in \mathbb{B}$. Fixing a coefficient $\rho \in \mathcal{D}$, we now consider for every index $j \in \mathbb{N}$ the j -th smallest eigenvalue $\lambda_j(\rho, \mathbf{k})$ as a function of \mathbf{k} from \mathbb{B} into \mathbb{R} .

Before we can prove our next result, we need to define some operators. Given a coefficient $\rho \in \mathcal{D}$ and a vector $\mathbf{k} \in \mathbb{B}$, we define the operators $\tilde{A}_{\mathbf{k}}(\rho) : \mathbf{W} \rightarrow \mathbf{W}^*$, $B_{\mathbf{k}} : \mathbf{Z} \rightarrow Q^*$, and $-\Delta_{\mathbf{k}} : Q \rightarrow Q^*$ by

$$\tilde{A}_{\mathbf{k}}(\rho)\mathbf{w}[\mathbf{v}] := a_{\mathbf{k}}(\rho)(\mathbf{w}, \mathbf{v}) \quad \text{for all } \mathbf{w}, \mathbf{v} \in \mathbf{W} \quad (4.69)$$

$$B_{\mathbf{k}}\mathbf{f}[q] := b_{\mathbf{k}}(\mathbf{f}, q) \quad \text{for all } \mathbf{f} \in \mathbf{Z}, q \in Q \quad (4.70)$$

$$-\Delta_{\mathbf{k}}p[q] := \int_{\Omega} \overline{(\nabla + i\mathbf{k})p} \cdot (\nabla + i\mathbf{k})q \quad \text{for all } p, q \in Q. \quad (4.71)$$

Recall that \mathbf{W} denotes the complex Hilbert space $\mathbf{H}_{\text{per}}(\text{curl}; \Omega)$, and that Q denotes the complex Hilbert space $H_{\text{per}}^1(\Omega)$. Notice that the operators $A_{\mathbf{k}}(\rho)$, defined by (4.43), differ from the operators $\tilde{A}_{\mathbf{k}}(\rho)$ only in the domain of definition and in the range.

Given a vector $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$, it is not difficult to show that for every continuous linear functional $f \in Q^*$ there exists a unique function $p \in Q$, such that

$$-\Delta_{\mathbf{k}}p = f. \quad (4.72)$$

One can also show that the solution p is bounded in Q and that it depends continuously on f . For $\mathbf{k} = \mathbf{0}$ the solution of (4.72) is only unique up to an additive constant. However, it can be chosen uniquely in the Hilbert space $Q_{\diamond} := Q \cap L_{\diamond}^2(\Omega)$. Recall that the Hilbert space $L_{\diamond}^2(\Omega)$ was defined in (4.24). It follows that for every $\mathbf{k} \in \mathbb{B}$ there exists a continuous linear operator $(-\Delta_{\mathbf{k}})^{-1} : Q^* \rightarrow Q$, which maps every functional $f \in Q^*$ either to the unique solution $p \in Q$ of (4.72), if $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$, or to the unique solution $p \in Q_{\diamond}$ of (4.72), if $\mathbf{k} = \mathbf{0}$. It follows that this operator $(-\Delta_{\mathbf{k}})^{-1}$ is the continuous, conjugate-linear inverse of $-\Delta_{\mathbf{k}}$.

Finally, we define for every vector $\mathbf{k} \in \mathbb{B}$ the operator $\pi_{\mathbf{k}} : \mathbf{W} \rightarrow \mathbf{W}$ by

$$\pi_{\mathbf{k}} := \text{id}_{\mathbf{W}} - (\nabla + i\mathbf{k})(-\Delta_{\mathbf{k}})^{-1}B_{\mathbf{k}}|_{\mathbf{W}}, \quad (4.73)$$

where $\text{id}_{\mathbf{W}}$ denotes the identity on \mathbf{W} . The utility of the operator $\pi_{\mathbf{k}}$ is made obvious by the following proposition.

Proposition 4.23. *For every $\mathbf{k} \in \mathbb{B}$, the operator $\pi_{\mathbf{k}}$ is the \mathbf{Z} -orthogonal projection from \mathbf{W} onto $\mathbf{V}_{\mathbf{k}}$.*

Proof. Let $\mathbf{w} \in \mathbf{W}$ be an arbitrary function. By Theorem 4.8(a) there exist for every $\mathbf{k} \in \mathbb{B} \setminus \{\mathbf{0}\}$ uniquely determined functions $\boldsymbol{\psi} \in H_{\text{per}}^1(\Omega)^3$ and $p \in Q$, such that $\mathbf{w} = (\nabla + i\mathbf{k}) \times \boldsymbol{\psi} + (\nabla + i\mathbf{k})p$. According to the definition of the operators $B_{\mathbf{k}}$ and $(-\Delta_{\mathbf{k}})^{-1}$, we have $(-\Delta_{\mathbf{k}})^{-1}B_{\mathbf{k}}\mathbf{w} = p$ and hence $\pi_{\mathbf{k}}\mathbf{w} = (\nabla + i\mathbf{k}) \times \boldsymbol{\psi}$. By Proposition 4.11 we have that $\pi_{\mathbf{k}}\mathbf{w} \in \mathbf{V}_{\mathbf{k}}$. Furthermore, since $\mathbf{w} - \pi_{\mathbf{k}}\mathbf{w} = (\nabla + i\mathbf{k})p$, we have by Corollary 4.9(a) that $\pi_{\mathbf{k}}$ is indeed an \mathbf{Z} -orthogonal projection and onto.

Next, we consider the case $\mathbf{k} = \mathbf{0}$. By Theorem 4.8(b) there exist uniquely determined function $\boldsymbol{\psi} \in H_{\text{per},\diamond}^1(\Omega)^3$ and $p \in Q_{\diamond}$, as well as a uniquely determined vector $\mathbf{z} \in \mathbb{C}^3$, such that $\mathbf{w} = \nabla \times \boldsymbol{\psi} + \nabla p + \mathbf{z}$. We find $(-\Delta_{\mathbf{0}})^{-1}B_{\mathbf{0}}\mathbf{w} = \nabla p$, which implies that $\pi_{\mathbf{0}}\mathbf{w} = \nabla \times \boldsymbol{\psi} + \mathbf{z}$. Again, we have by Proposition 4.11 that $\pi_{\mathbf{0}}\mathbf{w} \in \mathbf{V}_{\mathbf{0}}$, and by Corollary 4.9(b) that $\pi_{\mathbf{0}}$ is an \mathbf{Z} -orthogonal projection and onto. \square

With the operators $\tilde{A}_{\mathbf{k}}(\rho)$ and $\pi_{\mathbf{k}}$ at hand, we are now able to prove the following regularity result for the functions $\mathbf{k} \mapsto \lambda_j(\rho, \mathbf{k})$.

Proposition 4.24. *For every index $j \in \mathbb{N}$ and every coefficient $\rho \in \mathcal{D}$ the functional $\lambda_j(\rho, \cdot) : \mathbb{B} \rightarrow \mathbb{R}$ is continuous.*

Proof. First, we notice that $\tilde{A}_{\mathbf{k}}(\rho) = \pi_{\mathbf{k}}^* A_{\mathbf{k}}(\rho) \pi_{\mathbf{k}}$ for all $\rho \in \mathcal{D}$ and all $\mathbf{k} \in \mathbb{B}$, where $\pi_{\mathbf{k}}^*$ denotes the dual operator of $\pi_{\mathbf{k}}$. From this we deduce that the positive eigenvalues of the operators $\tilde{A}_{\mathbf{k}}(\rho)$ and $A_{\mathbf{k}}(\rho)$ coincide for every $\rho \in \mathcal{D}$ and every $\mathbf{k} \in \mathbb{B}$. Choosing a fixed coefficient $\rho \in \mathcal{D}$ and an arbitrary vector $\mathbf{k} \in \mathbb{B}$, we have

$$\begin{aligned} \tilde{A}_{\mathbf{k}}(\rho)\mathbf{w}[\mathbf{v}] &= \int_{\Omega} \rho \left(\overline{\nabla \times \mathbf{w}} \cdot \nabla \times \mathbf{v} + \overline{\nabla \times \mathbf{w}} \cdot [\mathbf{i}\mathbf{k}]_{\times} \mathbf{v} \right. \\ &\quad \left. - [\mathbf{i}\mathbf{k}]_{\times} \overline{\mathbf{w}} \cdot \nabla \times \mathbf{v} - [\mathbf{i}\mathbf{k}]_{\times} \overline{\mathbf{w}} \cdot [\mathbf{i}\mathbf{k}]_{\times} \mathbf{v} \right) \\ &= \int_{\Omega} \rho \left(\overline{\nabla \times \mathbf{w}} \cdot \nabla \times \mathbf{v} + [\mathbf{i}\mathbf{k}]_{\times}^* \overline{\nabla \times \mathbf{w}} \cdot \mathbf{v} \right. \\ &\quad \left. - [\mathbf{i}\mathbf{k}]_{\times} \overline{\mathbf{w}} \cdot \nabla \times \mathbf{v} - [\mathbf{i}\mathbf{k}]_{\times}^* [\mathbf{i}\mathbf{k}]_{\times} \overline{\mathbf{w}} \cdot \mathbf{v} \right) \quad \text{for all } \mathbf{w}, \mathbf{v} \in \mathbf{W}, \end{aligned}$$

where $[\mathbf{i}\mathbf{k}]_{\times}$ denotes the cross product matrix corresponding to $i\mathbf{k}$ (see Section 2.2). Now, let $\{\mathbf{k}^{(n)}\}_{n \in \mathbb{N}}$ be a sequence in \mathbb{B} , which converges to \mathbf{k} . Choosing an arbitrary number $n \in \mathbb{N}$ we obtain the estimate

$$|(\tilde{A}_{\mathbf{k}^{(n)}}(\rho) - \tilde{A}_{\mathbf{k}}(\rho))\mathbf{w}[\mathbf{v}]| \leq \|\rho\|_{\infty} (2f_n + g_n) \|\mathbf{w}\|_{\text{curl},\Omega} \|\mathbf{v}\|_{\text{curl},\Omega} \quad \text{for all } \mathbf{w}, \mathbf{v} \in \mathbf{W},$$

where, the numbers f_n and g_n are given by

$$\begin{aligned} f_n &:= \|\mathbf{k}^{(n)} - \mathbf{k}\|_{\times} && \text{for all } n \in \mathbb{N}, \\ g_n &:= \|\mathbf{k}^{(n)*}[\mathbf{k}^{(n)}]_{\times}^* - [\mathbf{k}]_{\times}^*[\mathbf{k}]_{\times}^*\|_2 && \text{for all } n \in \mathbb{N}. \end{aligned}$$

Here, we denote by $\|\cdot\|_2$ the spectral norm on $\mathbb{C}^{3 \times 3}$. Clearly, we have $f_n \rightarrow 0$ and $g_n \rightarrow 0$ as $n \rightarrow \infty$. The above estimate hence implies

$$\|\tilde{A}_{\mathbf{k}^{(n)}}(\rho) - \tilde{A}_{\mathbf{k}}(\rho)\|_{\mathbf{W}, \mathbf{W}^*} \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

We have thus shown that the operator $\tilde{A}_{\mathbf{k}^{(n)}}(\rho)$ converges to $\tilde{A}_{\mathbf{k}}(\rho)$ as $n \rightarrow \infty$. This convergence in particular implies the generalized convergence of $\tilde{A}_{\mathbf{k}^{(n)}}(\rho)$ to $\tilde{A}_{\mathbf{k}}(\rho)$ in the sense of Kato, since the assumption of Theorem 2.24 in Chapter IV, Section 2.6 in [45] hold a fortiori. Kato establishes in Chapter IV, Section 3.5 in [45] that the generalized convergence implies the convergence of the operators' spectra. Therefore, we have $\lambda_j(\rho, \mathbf{k}^{(n)}) \rightarrow \lambda_j(\rho, \mathbf{k})$ as $n \rightarrow \infty$, which completes this proof. \square

Recall that the unknown eigenvalues λ in Problem 4.16 represent frequencies ω of the unknown Bloch modes of a photonic crystal via $\lambda = \omega^2/c_0^2$ (see Section 4.1). Assuming $c_0 = 1$ in the following, we defined for every coefficient $\rho \in \mathcal{D}$ and every index $j \in \mathbb{N}$ the function $\omega_j : \mathcal{D} \times \mathbb{B} \rightarrow \mathbb{R}$ by

$$\omega_j(\rho, \mathbf{k}) := \sqrt{\lambda_j(\rho, \mathbf{k})} \quad \text{for all } \rho \in \mathcal{D}, \mathbf{k} \in \mathbb{B}. \quad (4.74)$$

By Proposition 4.24 the functions $\omega_j(\rho, \cdot) : \mathbb{B} \rightarrow \mathbb{R}$ are continuous for every every index $j \in \mathbb{N}$ and every coefficient $\rho \in \mathcal{D}$. Provided that ρ represents the medium structure of a photonic crystal via $\rho = 1/\varepsilon_r$, where ε_r denotes the crystal's relative electric permittivity field, the graphs of the functions $\omega_1(\rho, \cdot), \omega_2(\rho, \cdot), \dots$, etc. are referred to as the *photonic bands* of that crystal. In particular, the graph of the function $\omega_j(\rho, \cdot)$ is called the *j -th photonic band* for $i \in \mathbb{N}$. The union of all photonic bands is called the *photonic band structure*.

Given an index $j \in \mathbb{N}$ we shall say that the j -th and the $(j+1)$ -th photonic band are *strictly separated*, if

$$\omega_j(\rho, \mathbf{k}) < \omega_{j+1}(\rho, \mathbf{k}) \quad \text{for all } \mathbf{k} \in \mathbb{B}. \quad (4.75)$$

We shall say that there exists a *band gap* between the j -th and the $(j+1)$ -th photonic band, if

$$\max_{\mathbf{k} \in \mathbb{B}} \omega_j(\rho, \mathbf{k}) < \min_{\mathbf{k} \in \mathbb{B}} \omega_{j+1}(\rho, \mathbf{k}). \quad (4.76)$$

The corresponding *gap width* $w_j(\rho)$ is given by

$$w_j(\rho) := \min_{\mathbf{k} \in \mathbb{B}} \omega_{j+1}(\rho, \mathbf{k}) - \max_{\mathbf{k} \in \mathbb{B}} \omega_j(\rho, \mathbf{k}). \quad (4.77)$$

In Section 3.2 we discussed that the complete symmetry of a photonic crystal's medium structure consists of translational symmetries as well as of non-translational symmetries. The translational symmetries are captured by the crystal's Bravais lattice Λ . The non-translational symmetries are characterized in terms of so-called crystallographic point groups. Every element of a photonic crystal's point group corresponds to an orthogonal transformation that constitutes a symmetry operation for the photonic crystal. The following proposition reveals, how the non-translational symmetries of a three-dimensional photonic crystal's medium structure influence the corresponding band structure.

Proposition 4.25. *Let $\boldsymbol{\theta} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ be an orthogonal transformation, which constitutes a symmetry operation for a given coefficient $\rho \in \mathcal{D}$, as well as for the first Brillouin zone \mathbb{B} . Then, we have*

$$\begin{aligned} \lambda_j(\rho, \boldsymbol{\theta}^{-1}(\mathbf{k})) &= \lambda_j(\rho, \mathbf{k}) && \text{for all } j \in \mathbb{N}, \mathbf{k} \in \mathbb{B}, \\ \mathbf{u}_j(\rho, \boldsymbol{\theta}^{-1}(\mathbf{k})) &= \boldsymbol{\theta}^{-1}(\mathbf{u}_j(\rho, \mathbf{k}) \circ \boldsymbol{\theta}) && \text{for all } j \in \mathbb{N}, \mathbf{k} \in \mathbb{B}. \end{aligned}$$

Proof. To simplify notations we define for given $j \in \mathbb{N}$, $\rho \in \mathcal{D}$, and $\mathbf{k} \in \mathbb{B}$ the real number $\lambda := \lambda_j(\rho, \mathbf{k})$, and the vector field $\mathbf{u} := \mathbf{u}_j(\rho, \mathbf{k})$. Furthermore, let $\boldsymbol{\theta}$ be given by $\boldsymbol{\theta}(\mathbf{x}) := \boldsymbol{\Theta}\mathbf{x}$ for all $\mathbf{x} \in \mathbb{R}^3$, where $\boldsymbol{\Theta} \in \text{O}_3(\mathbb{R})$ is an orthogonal matrix. For every diffeomorphism $\boldsymbol{\varphi} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$, which satisfies $\boldsymbol{\varphi}(\Omega) = \Omega$, we define the mapping $(\cdot)_{\boldsymbol{\varphi}} : \mathbf{Z} \rightarrow \mathbf{Z}$ by

$$\mathbf{f}_{\boldsymbol{\varphi}} := (\boldsymbol{\varphi}')^{\text{T}}(\mathbf{f} \circ \boldsymbol{\varphi}) \quad \text{for all } \mathbf{f} \in \mathbf{Z}.$$

Clearly, each such mapping $(\cdot)_{\boldsymbol{\varphi}}$ is a linear isomorphism, where the inverse mapping of $(\cdot)_{\boldsymbol{\varphi}}$ is given by $(\cdot)_{\boldsymbol{\varphi}^{-1}}$. Moreover, it can be shown (see e.g. Lemma 8 in Section A.1 in [25]) that

$$\nabla \times \mathbf{w}_{\boldsymbol{\varphi}} = \det(\boldsymbol{\varphi}') (\boldsymbol{\varphi}')^{-1} (\nabla \times \mathbf{w} \circ \boldsymbol{\varphi}) \quad \text{for all } \mathbf{w} \in \mathbf{W}.$$

One observes that $\nabla \times \mathbf{w}_{\boldsymbol{\theta}}$ is given by the so-called *Piola transform* of $\nabla \times \mathbf{w}$ with respect to $\boldsymbol{\varphi}$ (see Definition 7.18 in [53]). Accordingly, we obtain $\mathbf{f}_{\boldsymbol{\theta}} = \boldsymbol{\Theta}^{\text{T}}(\mathbf{f} \circ \boldsymbol{\theta})$ for all $\mathbf{f} \in \mathbf{Z}$, as well as $\nabla \times \mathbf{w}_{\boldsymbol{\theta}} = \det(\boldsymbol{\Theta})\boldsymbol{\Theta}^{\text{T}}(\nabla \times \mathbf{w} \circ \boldsymbol{\theta})$ for all $\mathbf{w} \in \mathbf{W}$. It follows that the mapping $(\cdot)_{\boldsymbol{\theta}}$ constitutes an isometric, linear isomorphism from \mathbf{Z} onto \mathbf{Z} , as well as from \mathbf{W} onto \mathbf{W} . We now show that the mapping $(\cdot)_{\boldsymbol{\theta}}$ also constitutes an isometric linear isomorphism from $\mathbf{V}_{\mathbf{k}}$ onto $\mathbf{V}_{\boldsymbol{\theta}^{-1}(\mathbf{k})}$. To this end, we define for every diffeomorphism $\boldsymbol{\varphi} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$, which satisfies $\boldsymbol{\varphi}(\Omega) = \Omega$, the

mapping $(\cdot)^\varphi : Q \rightarrow Q$ by $q^\varphi := q \circ \varphi$ for all $q \in Q$. Since $\nabla q^\theta = \Theta^T((\nabla q) \circ \theta)$ for all $q \in Q$, one easily shows that the mapping $(\cdot)^\theta$ is an isometrical, linear isomorphism. Clearly, the inverse mapping of $(\cdot)^\theta$ is given by $(\cdot)^{\theta^{-1}}$. Choosing $\mathbf{v} \in \mathbf{V}_k$ and $q \in Q$ arbitrarily, and letting $r := q^{\theta^{-1}}$, we find that

$$\begin{aligned}
b_{\theta^{-1}(\mathbf{k})}(\mathbf{v}_\theta, q) &= b_{\theta^{-1}(\mathbf{k})}(\mathbf{v}_\theta, r^\theta) = \int_{\Omega} \Theta^T(\bar{\mathbf{v}} \circ \theta) \cdot (\nabla + i\theta^{-1}(\mathbf{k}))(r \circ \theta) \\
&= \int_{\Omega} \Theta^T(\bar{\mathbf{v}} \circ \theta) \cdot \Theta^T((\nabla r \circ \theta) + i\mathbf{k}(r \circ \theta)) \\
&= \int_{\Omega} (\bar{\mathbf{v}} \circ \theta) \cdot ((\nabla + i\mathbf{k})r \circ \theta) = \int_{\theta(\Omega)} \bar{\mathbf{v}} \cdot (\nabla + i\mathbf{k})r \\
&= b_{\mathbf{k}}(\mathbf{v}, r) \\
&= 0.
\end{aligned}$$

Since this identity can be established for all $q \in Q$, we have that $\mathbf{v}_\theta \in \mathbf{V}_{\theta^{-1}(\mathbf{k})}$. It follows that the image of \mathbf{V}_k under the mapping $(\cdot)_\theta$ is a subset of $\mathbf{V}_{\theta^{-1}(\mathbf{k})}$. By a similar computation one can show, that for each function $\mathbf{v} \in \mathbf{V}_{\theta^{-1}(\mathbf{k})}$ we have $\mathbf{v}_{\theta^{-1}} \in \mathbf{V}_k$, implying that the restriction of $(\cdot)_\theta$ to \mathbf{V}_k is surjective. Injectivity, and hence bijectivity, of this restriction follow from the fact that the mapping $(\cdot)_\theta$ is an isometry from \mathbf{W} onto \mathbf{W} .

Since Θ is an orthogonal matrix, we have $(\Theta \mathbf{x}) \times (\Theta \mathbf{y}) = \det(\Theta) \Theta(\mathbf{x} \times \mathbf{y})$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{C}^3$ according to (2.6) and (2.8). Choosing $\mathbf{v} \in \mathbf{V}_k$ arbitrarily, we obtain

$$\begin{aligned}
a_{\theta^{-1}(\mathbf{k})}(\rho)(\mathbf{u}_\theta, \mathbf{v}_\theta) &= \int_{\Omega} \overline{\rho(\nabla + i\theta^{-1}(\mathbf{k})) \times \mathbf{u}_\theta} \cdot (\nabla + i\theta^{-1}(\mathbf{k})) \times \mathbf{v}_\theta \\
&= \int_{\Omega} \overline{\rho(\nabla \times \mathbf{u}_\theta + i\Theta^T \mathbf{k} \times \mathbf{u}_\theta)} \cdot (\nabla \times \mathbf{v}_\theta + i\Theta^T \mathbf{k} \times \mathbf{v}_\theta) \\
&= \int_{\Omega} \overline{\rho((\nabla + i\mathbf{k}) \times \mathbf{u} \circ \theta)} \cdot ((\nabla + i\mathbf{k}) \times \mathbf{v} \circ \theta) \\
&= \int_{\theta(\Omega)} (\rho \circ \theta^{-1}) \overline{(\nabla + i\mathbf{k}) \times \mathbf{u}} \cdot (\nabla + i\mathbf{k}) \times \mathbf{v} \\
&= a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{v}).
\end{aligned}$$

Similarly, we find that $m(\mathbf{u}_\theta, \mathbf{v}_\theta) = m(\mathbf{u}, \mathbf{v})$. Since both identities can be established for all $\mathbf{v} \in \mathbf{V}_k$, we have that $a_{\theta^{-1}(\mathbf{k})}(\rho)(\mathbf{u}_\theta, \mathbf{v}_\theta) - \lambda m(\mathbf{u}_\theta, \mathbf{v}_\theta) = 0$ for all $\mathbf{v} \in \mathbf{V}_k$. Since $(\cdot)_\theta$ is a linear isomorphism from \mathbf{V}_k onto $\mathbf{V}_{\theta^{-1}(\mathbf{k})}$, we also have

$$a_{\theta^{-1}(\mathbf{k})}(\rho)(\mathbf{u}_\theta, \mathbf{v}) = \lambda m(\mathbf{u}_\theta, \mathbf{v}) \quad \text{for all } \mathbf{v} \in \mathbf{V}_{\theta^{-1}(\mathbf{k})},$$

implying that $\lambda_j(\rho, \theta^{-1}(\mathbf{k})) = \lambda$, and $\mathbf{u}_j(\rho, \theta^{-1}(\mathbf{k})) = \mathbf{u}_\theta = \Theta^T(\mathbf{u}_j(\rho, \mathbf{k}) \circ \theta)$. \square

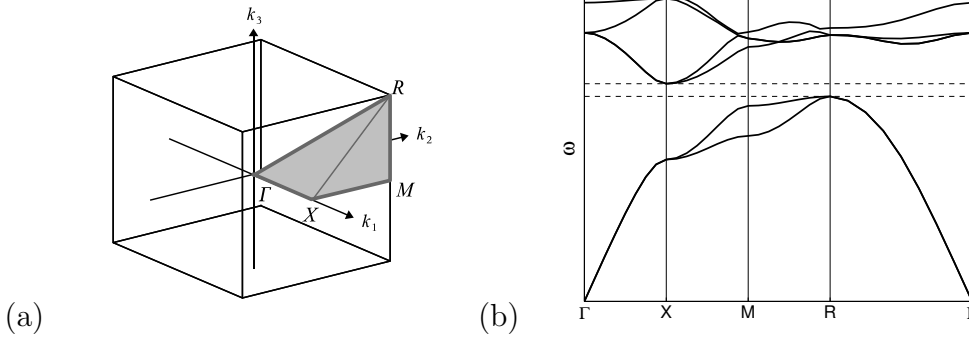


Figure 4.1: The first Brillouin zone with the irreducible zone (dark) and the critical points (a). The band diagram (b) shows the band structure along the critical path Γ - X - M - R - Γ . The dashed lines indicate a suspected band gap. Here, the space group is isomorphic to the semidirect product $\mathbb{Z}^3 \rtimes O_3(\mathbb{Z})$.

In the following we illustrate a practical consequence of Proposition 4.25. Suppose that ε_r is the relative electric permittivity function of a three-dimensional photonic crystal, whose medium structure is periodic with respect to a simple cubic lattice $\Lambda = a\mathbb{Z}^3$. Here $a > 0$ denotes the lattice constant. Furthermore, suppose that the crystal's point group is isomorphic to $O_3(\mathbb{Z})$. Letting

$$G := \{ \boldsymbol{\theta} : \mathbb{R}^3 \rightarrow \mathbb{R}^3 \mid \boldsymbol{\theta}(\mathbf{x}) = \boldsymbol{\Theta}\mathbf{x}, \boldsymbol{\Theta} \in O_3(\mathbb{Z}) \},$$

we then have that

$$\varepsilon_r \circ \boldsymbol{\theta} = \varepsilon_r \quad \text{for all } \boldsymbol{\theta} \in G.$$

The crystal's Wigner-Seitz cell is given by the cube $[-a/2, a/2]^3$. Letting $\Omega := (-a/2, a/2)^3$ and $\rho := 1/\varepsilon_r|_{\Omega}$, we also have that

$$\rho \circ \boldsymbol{\theta} = \rho \quad \text{for all } \boldsymbol{\theta} \in G.$$

One easily verifies that the reciprocal lattice of Λ is also a simple cubic lattice, which is given by $\widehat{\Lambda} = (2\pi/a)\mathbb{Z}^3$. It thus follows that the first Brillouin zone of Λ is given by $\mathbb{B} = [-\pi/a, \pi/a]^3$. Clearly, this first Brillouin zone is symmetric with respect to G . By Proposition 4.25 we thus have that

$$\lambda_j(\rho, \boldsymbol{\theta}^{-1}(\mathbf{k})) = \lambda_j(\rho, \mathbf{k}) \quad \text{for all } j \in \mathbb{N}, \mathbf{k} \in \mathbb{B}, \boldsymbol{\theta} \in G. \quad (4.78)$$

From (4.78) we deduce that for every index $j \in \mathbb{N}$ the function $\lambda_j(\rho, \cdot) : \mathbb{B} \rightarrow \mathbb{R}$ is G -symmetric, and hence uniquely determined by its restriction to any fundamental region of G . The closure of such a fundamental region is called an *irreducible zone*.

One can show that the tetrahedron

$$K := \overline{\text{conv}}\{\mathbf{k}_\Gamma, \mathbf{k}_X, \mathbf{k}_M, \mathbf{k}_R\},$$

where

$$\mathbf{k}_\Gamma := \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}, \quad \mathbf{k}_X := \begin{pmatrix} \pi/a \\ 0 \\ 0 \end{pmatrix}, \quad \mathbf{k}_M := \begin{pmatrix} \pi/a \\ \pi/a \\ 0 \end{pmatrix}, \quad \mathbf{k}_R := \begin{pmatrix} \pi/a \\ \pi/a \\ \pi/a \end{pmatrix},$$

is an irreducible zone for every simple cubic photonic crystal with lattice constant $a > 0$, whose crystallographic point group is isomorphic to $O_3(\mathbb{Z})$. The vectors \mathbf{k}_Γ , \mathbf{k}_X , \mathbf{k}_M and \mathbf{k}_R are referred to as *critical points*. The path line segments, which connect the critical points \mathbf{k}_Γ , \mathbf{k}_X , \mathbf{k}_M , \mathbf{k}_R and \mathbf{k}_X (in this order), is called the *critical path*. We depicted the tetrahedron K , the critical point, and the critical path in Figure 4.1(a). By plotting the graphs of the functions $\omega_1(\rho, \cdot), \dots, \omega_n(\rho, \cdot)$ along the critical path for some $n \in \mathbb{N}$, one obtains a so-called *photonic band diagram*. An example of such a band diagram is depicted in Figure 4.1(b).

In the literature on photonic crystals it is common practice to assume that the functions $\omega_j(\rho, \cdot)$ attain their extrema exactly on the critical path. Provided that this assumption indeed holds true, a band gap is revealed by a region in the band diagram, which is bounded by two horizontal lines such that neither of the two lines is crossed by any of the curves in the band diagram. The distance between the two horizontal lines then corresponds to the gap width.

4.7 The Two-Dimensional Case

So far, we only considered the spectral theory for the weakly formulated eigenvalue problems that determine the band structures of three-dimensional photonic crystals. In this section we turn to the eigenvalue problems that arise for TM- and TE-polarized Bloch modes in two-dimensional photonic crystals (see Section 3.5). The spectral theory related to these eigenvalue problems has already been discussed thoroughly in the literature (see e.g. [29], [50]), which is why we only give a brief account of the most important results.

Throughout this section we assume that $\Lambda \subset \mathbb{R}^2$ is a Bravais lattice of rank 2, and that $\Omega \subset \mathbb{R}^2$ is a primitive domain of Λ . As before, we denote by \mathbb{B} the first Brillouin zone of Λ . For convenience we define the complex Hilbert spaces

$$V := H_{\text{per}}^1(\Omega), \tag{4.79}$$

$$Z := L^2(\Omega). \tag{4.80}$$

Furthermore, we assume that the real Banach space \mathcal{E} and the open function set \mathcal{D} are defined according to (4.31) and (4.32).

Suppose that the Bravais lattice Λ represents the periodicity a two-dimensional photonic crystal's medium structure. In Section 3.5 we derived an eigenvalue problem, which determines the existence of so-called TM-polarized Bloch modes with a given quasimomentum vector $\mathbf{k} \in \mathbb{B}$ (see (3.49) on page 43). In order to establish a weak formulation of this eigenvalue problem, we define for every coefficient $\rho \in \mathcal{D}$ and every vector $\mathbf{k} \in \mathbb{B}$ the sesquilinear forms $a_{\mathbf{k}}^{\text{TM}} : V \times V \rightarrow \mathbb{C}$ and $m^{\text{TM}}(\rho) : V \times V \rightarrow \mathbb{C}$ by

$$a_{\mathbf{k}}^{\text{TM}}(u, v) := \int_{\Omega} \overline{(\nabla + \mathbf{i}\mathbf{k})u} \cdot (\nabla + \mathbf{i}\mathbf{k})v \quad \text{for all } u, v \in V, \quad (4.81)$$

$$m^{\text{TM}}(\rho)(f, g) := \int_{\Omega} \rho \bar{f}g \quad \text{for all } f, g \in Z. \quad (4.82)$$

With this, we consider the following problem.

Problem 4.26. *Given a coefficient $\rho \in \mathcal{D}$ and a vector $\mathbf{k} \in \mathbb{B}$, find eigenvalues $\lambda \in \mathbb{C}$ and corresponding eigenfunctions $u \in V$, such that*

$$a_{\mathbf{k}}^{\text{TM}}(u, v) = \lambda m^{\text{TM}}(\rho)(u, v) \quad \text{for all } v \in V. \quad (4.83)$$

If $\rho = \varepsilon_r|_{\Omega}$, where ε_r is the relative electric permittivity function of a two-dimensional photonic crystal, the (4.83) is precisely the weak formulation of the eigenvalue equation (3.49) on page 43, and thus determines the existence of TM-polarized Bloch modes in the photonic crystal. One can show that for every vector $\mathbf{k} \in \mathbb{B}$ the sesquilinear form $a_{\mathbf{k}}(\rho)$ satisfies

$$\begin{aligned} a_{\mathbf{k}}^{\text{TM}}(v, u) &= \overline{a_{\mathbf{k}}^{\text{TM}}(u, v)} && \text{for all } v, u \in V, \\ |a_{\mathbf{k}}^{\text{TM}}(v, u)| &\leq \beta^{\text{TM}}(\mathbf{k}) \|v\|_{1, \Omega} \|u\|_{1, \Omega} && \text{for all } u, v \in V, \\ a_{\mathbf{k}}^{\text{TM}}(v, v) &\geq \alpha^{\text{TM}}(\mathbf{k}) \|v\|_{1, \Omega}^2 - \kappa^{\text{TM}}(\mathbf{k}) \|v\|_{\Omega}^2 && \text{for all } v \in V, \\ a_{\mathbf{k}}^{\text{TM}}(v, v) &\geq 0 && \text{for all } v \in V. \end{aligned}$$

The functions $\alpha^{\text{TM}}, \beta^{\text{TM}}, \kappa^{\text{TM}} : \mathbb{B} \rightarrow \mathbb{R}$ are given by

$$\alpha^{\text{TM}}(\mathbf{k}) := \left(1 - \frac{|\mathbf{k}|}{\text{diam}(\mathbb{B})}\right) \quad \text{for all } \mathbf{k} \in \mathbb{B}, \quad (4.84)$$

$$\beta^{\text{TM}}(\mathbf{k}) := 2 \max\{1, |\mathbf{k}|^2\} \quad \text{for all } \mathbf{k} \in \mathbb{B}, \quad (4.85)$$

$$\kappa^{\text{TM}}(\mathbf{k}) := \alpha^{\text{TM}}(\mathbf{k}) + (\text{diam}(\mathbb{B})|\mathbf{k}| + |\mathbf{k}|^2) \quad \text{for all } \mathbf{k} \in \mathbb{B}. \quad (4.86)$$

Thus, it follows that the sesquilinear form $a_{\mathbf{k}}^{\text{TM}}$ is conjugate-symmetric, bounded, V - Z -coercive and positive semidefinite for every $\mathbf{k} \in \mathbb{B}$.

By a line of argument similar to the one presented in Section 4.4, one can show that for every coefficient $\rho \in \mathcal{D}$ and every vector $\mathbf{k} \in \mathbb{B}$ there exists an increasing sequence $\{\lambda_j^{\text{TM}}(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$ of non-negative, real eigenvalues of the weak eigenvalue equation (4.83). Moreover, one can show that the sequence of the corresponding eigenfunctions $\{u_j^{\text{TM}}(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$ forms a complete system of V , which is orthogonal with respect to the inner product $\langle \cdot, \cdot \rangle_\rho$ on Z . For every $\rho \in \mathcal{D}$ the inner product $\langle \cdot, \cdot \rangle_\rho : Z \times Z \rightarrow \mathbb{C}$ is given by

$$\langle f, g \rangle_\rho := m^{\text{TM}}(\rho)(f, g) \quad \text{for all } f, g \in Z. \quad (4.87)$$

By $\|\cdot\|_\rho$, we denote the corresponding norm on V . Clearly, for every $\rho \in \mathcal{D}$, one has that

$$\text{ess inf}_\Omega(\rho)\|f\|_\Omega \leq \|f\|_\rho \leq \text{ess sup}_\Omega(\rho)\|f\|_\Omega \quad \text{for all } f \in Z,$$

i.e., the norms $\|\cdot\|_\Omega$ and $\|\cdot\|_\rho$ are equivalent.

As in the three-dimensional setting, the eigenspace of every eigenvalue of (4.83) is finite-dimensional, and the j -th smallest eigenvalues $\lambda_j^{\text{TM}}(\rho, \mathbf{k})$ can be characterized by the Min-Max Principle as

$$\lambda_j^{\text{TM}}(\rho, \mathbf{k}) = \min_{\substack{U \subseteq V \\ \dim U = j}} \max_{u \in U \setminus \{0\}} \frac{a_{\mathbf{k}}^{\text{TM}}(u, u)}{m^{\text{TM}}(\rho)(u, u)} \quad \text{for all } j \in \mathbb{N} \quad (4.88)$$

or, alternatively, by Auchmuty's principle. In order to formulate the latter one, we define for every coefficient $\rho \in \mathcal{D}$ and every index $j \in \mathbb{N}$ the linear function space

$$V_{\mathbf{k},j}^{\text{TM}}(\rho) := \text{span}\{u_1^{\text{TM}}(\rho, \mathbf{k}), \dots, u_j^{\text{TM}}(\rho, \mathbf{k})\}, \quad (4.89)$$

as well as the projection operator $P_{\mathbf{k},j}^{\text{TM}}(\rho) : Z \rightarrow V_{\mathbf{k},j}(\rho)$ by

$$P_{\mathbf{k},j}^{\text{TM}}(\rho)f := \sum_{i=1}^j \frac{\langle u_i^{\text{TM}}(\rho, \mathbf{k}), f \rangle_\rho}{\|u_i^{\text{TM}}(\rho, \mathbf{k})\|_\rho^2} u_i^{\text{TM}}(\rho, \mathbf{k}) \quad \text{for all } f \in Z. \quad (4.90)$$

Then, Auchmuty's principle reads

$$\frac{-1}{2(\lambda_j^{\text{TM}}(\rho, \mathbf{k}) + \kappa^{\text{TM}}(\mathbf{k}))} = \min_{u \in V} \left(\frac{1}{2} (a_{\mathbf{k}}^{\text{TM}}(u, u) + \kappa^{\text{TM}}(\mathbf{k}) m^{\text{TM}}(\rho)(u, u)) - \|u - P_{\mathbf{k},j-1}^{\text{TM}}(\rho)u\|_\rho \right) \quad \text{for all } j \in \mathbb{N}. \quad (4.91)$$

In Section 3.5 we also derived an eigenvalue problem that determines the existence of so-called TE-polarized Bloch modes (see (3.50) on page 43) in two-dimensional photonic crystals. In order to establish also a weak formulation of this problem, we defined for every coefficient $\rho \in \mathcal{D}$ and every vector $\mathbf{k} \in \mathbb{B}$, we define the sesquilinear forms $a_{\mathbf{k}}^{\text{TE}}(\rho) : V \times V \rightarrow \mathbb{C}$ and $m^{\text{TE}} : Z \times Z \rightarrow \mathbb{C}$ by

$$a_{\mathbf{k}}^{\text{TE}}(\rho)(u, v) := \int_{\Omega} \rho \overline{(\nabla + \mathbf{i}\mathbf{k})u} \cdot (\nabla + \mathbf{i}\mathbf{k})v \quad \text{for all } u, v \in V, \quad (4.92)$$

$$m^{\text{TE}}(f, g) := \int_{\Omega} \bar{f}g \quad \text{for all } f, g \in Z. \quad (4.93)$$

We then consider the following problem.

Problem 4.27. *Given a coefficient $\rho \in \mathcal{D}$, and a vector $\mathbf{k} \in \mathbb{B}$, find eigenvalues $\lambda \in \mathbb{C}$ and corresponding eigenfunctions $u \in V$, such that*

$$a_{\mathbf{k}}^{\text{TE}}(\rho)(u, v) = \lambda m^{\text{TE}}(u, v) \quad \text{for all } v \in V. \quad (4.94)$$

One notices that (4.94) is precisely the weak formulation of the eigenvalue equation (3.49) on page 43, provided that the coefficient ρ is given by the relative electric permittivity function ε_r of a two-dimensional photonic crystal according to $\rho = 1/\varepsilon_r|_{\Omega}$. One can show that the sesquilinear form $a_{\mathbf{k}}^{\text{TE}}(\rho)$ is conjugate-symmetric, bounded, V - Z -coercive and positive semidefinite for every coefficient $\rho \in \mathcal{D}$ and every vector $\mathbf{k} \in \mathbb{B}$. It thus follows that for every $\rho \in \mathcal{D}$ and every $\mathbf{k} \in \mathbb{B}$ there exists an increasing sequence $\{\lambda_j^{\text{TE}}(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$ of non-negative, real eigenvalues of (4.94). The sequence of the corresponding eigenfunctions $\{u_j^{\text{TE}}(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$ forms a complete Z -orthogonal system of V .

In analogy to the definitions in Section 4.6, we define for every index $j \in \mathbb{N}$ the functions $\omega_j^{\text{TM}} : \mathcal{D} \times \mathbb{B} \rightarrow \mathbb{R}$ and $\omega_j^{\text{TE}} : \mathcal{D} \times \mathbb{B} \rightarrow \mathbb{R}$ by

$$\omega_j^{\text{TM}}(\rho, \mathbf{k}) := \sqrt{\lambda_j^{\text{TM}}(\rho, \mathbf{k})} \quad \text{for all } \rho \in \mathcal{D}, \mathbf{k} \in \mathbb{B}, \quad (4.95)$$

$$\omega_j^{\text{TE}}(\rho, \mathbf{k}) := \sqrt{\lambda_j^{\text{TE}}(\rho, \mathbf{k})} \quad \text{for all } \rho \in \mathcal{D}, \mathbf{k} \in \mathbb{B}. \quad (4.96)$$

Now suppose that ε_r is the relative electric permittivity function of a given, two-dimensional photonic crystal. Then, the graphs of the functions $\omega_1^{\text{TM}}(\varepsilon_r|_{\Omega}, \cdot)$, $\omega_2^{\text{TM}}(\varepsilon_r|_{\Omega}, \cdot), \dots$, etc. are referred to as the *TM bands* of the photonic crystal. Similarly, the graphs of the functions $\omega_1^{\text{TE}}(1/\varepsilon_r|_{\Omega}, \cdot), \omega_2^{\text{TE}}(1/\varepsilon_r|_{\Omega}, \cdot), \dots$, etc. are referred to as the *TE bands* of the photonic crystal. The union of all TM bands is called the *TM band structure*, and the union of all TE bands is called the *TE band structure* of the photonic crystal.

A two-dimensional photonic crystal with relative electric permittivity function ε_r is said to exhibit a *TM band gap*, if there exists an index $j \in \mathbb{N}$, such that

$$\max_{\mathbf{k} \in \mathbb{B}} \omega_j^{\text{TM}}(\varepsilon_r|_{\Omega}, \mathbf{k}) < \min_{\mathbf{k} \in \mathbb{B}} \omega_{j+1}^{\text{TM}}(\varepsilon_r|_{\Omega}, \mathbf{k})$$

The crystal is said to exhibit a *TE band gap*, if there exists an index $j \in \mathbb{N}$, such that

$$\max_{\mathbf{k} \in \mathbb{B}} \omega_j^{\text{TE}}(1/\varepsilon_r|_{\Omega}, \mathbf{k}) < \min_{\mathbf{k} \in \mathbb{B}} \omega_{j+1}^{\text{TE}}(1/\varepsilon_r|_{\Omega}, \mathbf{k}).$$

It should be noted that a TM band gap only corresponds to a range of inhibited TM wave frequencies, and that a TE band gap only corresponds to a range of inhibited TE wave frequencies. This means that a TM polarized wave with a frequency lying in a TE band gap of a two-dimensional photonic crystal might still be able to propagate inside the crystal. Therefore, two-dimensional photonic crystals only inhibit light propagation for frequencies, which belong to a so-called *complete band gap*. A crystal with a relative electric permittivity function ε_r is said to exhibit such a complete band gap, if there exist indices $i, j \in \mathbb{N}$, such that

$$\begin{aligned} & \max_{\mathbf{k} \in \mathbb{B}} \max \left\{ \omega_i^{\text{TM}}(\varepsilon_r|_{\Omega}, \mathbf{k}), \omega_j^{\text{TE}}(1/\varepsilon_r|_{\Omega}, \mathbf{k}) \right\} \\ & < \min_{\mathbf{k} \in \mathbb{B}} \min \left\{ \omega_{i+1}^{\text{TM}}(\varepsilon_r|_{\Omega}, \mathbf{k}), \omega_{j+1}^{\text{TE}}(1/\varepsilon_r|_{\Omega}, \mathbf{k}) \right\}. \end{aligned}$$

Clearly, a frequency ω_0 belongs to a complete band gap of a two-dimensional photonic crystal, if and only if ω_0 belongs to a TM band gap as well as to a TE band gap of that crystal.

Now, suppose that the relative electric permittivity function ε_r is periodic with respect to a square lattice $\Lambda = a\mathbb{Z}^2$ with lattice constant $a > 0$, and that the primitive domain Ω is given by the square $(-a/2, a/2)^2$. We further assume that the point group of ε_r is isomorphic to $O_2(\mathbb{Z})$. Then, we have by Lemma 4.1 in [29] that the functions $\omega_j^{\text{TM}}(\varepsilon_r|_{\Omega}, \cdot)$ and $\omega_j^{\text{TE}}(1/\varepsilon_r|_{\Omega}, \cdot)$ are uniquely determined by their restrictions to a fundamental region of the point group. As in the three-dimensional setting, the closure of such a fundamental region is called an irreducible zone.

For every two-dimensional photonic crystals, whose medium structure is periodic with respect to a square lattice with lattice constant $a > 0$, and whose point group is isomorphic to $O_2(\mathbb{Z})$, an irreducible zone is given by the triangle

$$K := \overline{\text{conv}}\{\mathbf{k}_{\Gamma}, \mathbf{k}_X, \mathbf{k}_M\},$$

where the critical points \mathbf{k}_{Γ} , \mathbf{k}_X , and \mathbf{k}_M are given by

$$\mathbf{k}_{\Gamma} := \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \mathbf{k}_X := \begin{pmatrix} \pi/a \\ 0 \end{pmatrix}, \quad \mathbf{k}_M := \begin{pmatrix} \pi/a \\ \pi/a \end{pmatrix}.$$

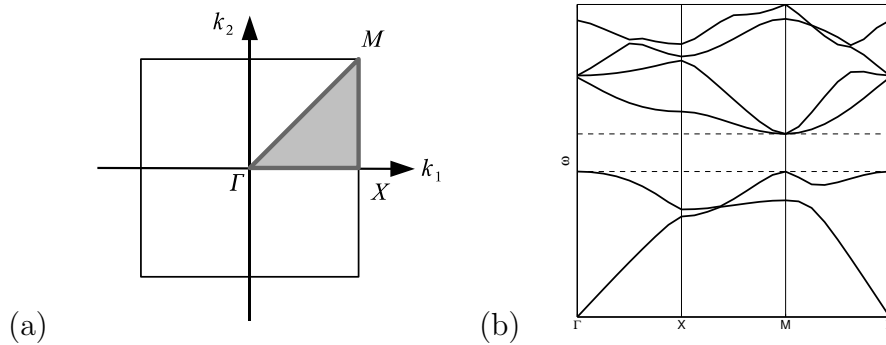


Figure 4.2: The first Brillouin zone with the irreducible zone (dark) and the critical points (a). The band diagram (b) shows the band structure along the critical path Γ - X - M - Γ . The dashed lines indicate a suspected band gap. Here, the space group is isomorphic to the semidirect product $\mathbb{Z}^2 \rtimes O_2(\mathbb{Z})$.

The corresponding critical path is defined as the boundary of K . In Figure 4.2(a) we depicted the first Brillouin zone \mathbb{B} and the irreducible zone K with its critical path and points. An example of a corresponding band diagram is depicted in Figure 4.2(b).

Chapter 5

Photonic Band Structure Optimization

In this chapter we study some mathematical aspects of photonic band structure optimization problems (PBSOPs). Recall that a PBSOP consists in finding medium structures of photonic crystals, which are optimal with respect to certain features of the crystals' photonic band structures. In Section 5.1 we establish a general framework, which allows us to treat a large class of PBSOPs as minimization problems over an admissible set of essentially bounded functions. In Section 5.2 we prove important results concerning the regularity of the goal functionals, which typically appear in the minimization problems. In particular, we show that these goal functionals are locally Lipschitz continuous and, under certain assumptions, even Lipschitz continuous. The existence of optimal solutions is discussed in the following two sections. In Section 5.3 we show that PBSOPs, which involve TM band structures of a two-dimensional photonic crystals, admit optimal solutions in an extended admissible set. In Section 5.4 we comment on why this result could not be shown so far for PBSOPs that involve TE band structures of two-dimensional photonic crystals or band structures of three-dimensional photonic crystals. Finally, we present some concrete optimization problems in Section 5.5.

5.1 The Formal Setting

With the definitions of Section 4.1, let $\mathcal{M}(\Omega)$ be the system of Lebesgue measurable subsets of Ω . Recall that Ω denotes a primitive cell of a Bravais lattice Λ of rank 3. Furthermore, let $\rho_{\min}, \rho_{\max} \in \mathbb{R}$ be two given constants, which satisfy

$$0 < \rho_{\min} < \rho_{\max}. \tag{5.1}$$

We define the set

$$\mathcal{C} := \left\{ \rho_{\min}(1 - \chi_S) + \rho_{\max}\chi_S \mid S \in \mathcal{M}(\Omega) \right\}, \quad (5.2)$$

where χ_S denotes the *characteristic function* of a set $S \subseteq \Omega$. Clearly, the set \mathcal{C} consists of all Lebesgue measurable functions on Ω with image in the set $\{\rho_0, \rho_1\}$. In particular, \mathcal{C} is a proper subset of the function set $\mathcal{D} \subset \mathcal{E}$, which we introduced in Section 4.3. In the following, we shall refer to the set \mathcal{C} as the *admissible set*.

In view of three-dimensional photonic crystals, \mathcal{C} represents the set of all possible Λ -periodic crystal structures (see Section 3.1). The constants ρ_{\min} and ρ_{\max} correspond to the reciprocal relative electric permittivities of two materials, the photonic crystal consists of. Given a function $\rho \in \mathcal{C}$, with $\rho = \rho_0(1 - \chi_S) + \rho_1\chi_S$ and $S \in \mathcal{M}(\Omega)$, the set S represents the region within the primitive domain Ω , which is occupied by the material whose reciprocal relative electric permittivity is given by ρ_{\max} . The remaining region $\Omega \setminus S$ is hence occupied by the material with reciprocal relative electric permittivity ρ_{\min} . The elements of the set \mathcal{C} will be referred to as *density functions*, henceforth.

As was discussed in Section 4.4, for every density $\rho \in \mathcal{C}$ and every vector $\mathbf{k} \in \mathbb{B}$ there exists an increasing sequence $\{\lambda_j(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$ of non-negative, real eigenvalues of the eigenvalue problem

$$(\nabla + i\mathbf{k}) \times [\rho(\nabla + i\mathbf{k}) \times \mathbf{u}] = \lambda \mathbf{u} \quad \text{in } \Omega$$

for $\lambda \in \mathbb{C}$ and $\mathbf{u} \in \mathbf{V}_{\mathbf{k}} \setminus \{\mathbf{0}\}$. Given an index $j \in \mathbb{N}$, the j -th smallest eigenvalue can be viewed as a function $\lambda_j : \mathcal{D} \times \mathbb{B} \rightarrow \mathbb{R}$. According to Proposition 4.24 the function $\lambda_j(\rho, \cdot) : \mathbb{B} \rightarrow \mathbb{R}$ is continuous for every density function $\rho \in \mathcal{C}$. The band structure of a photonic crystal, whose medium structure is represented by the density function $\rho \in \mathcal{C}$, is given by the graphs of all functions $\omega_1(\rho, \cdot)$, $\omega_2(\rho, \cdot)$, \dots , etc., where $\omega_j = \sqrt{\lambda_j}$ for every $j \in \mathbb{N}$.

Many optical properties of a photonic crystal, such as its effective refractive index or its transmission spectrum are determined by its band structure. The band structure itself is a function of the photonic crystal structure, which in our setting is given by a density function ρ belonging to the admissible set \mathcal{C} . Suppose that somebody wants to fabricate a three-dimensional photonic crystal with a specific optical property. Choosing two different materials, the question arises, how these two materials should be spatially arranged in order to obtain a photonic crystal with the desired optical property. In our setting the choice of the two materials determines the numbers ρ_{\min} and ρ_{\max} , and thus the admissible set \mathcal{C} . Devising a spatial arrangement of the two materials corresponds to choosing a specific element of \mathcal{C} . Thus, the problem is to find an element in \mathcal{C} , i.e. a density function, which is optimal with respect to the desired optical property.

In the following, we assume that the desired optical property of the photonic crystal can be expressed in terms of the lowest n bands of the crystal's band structure, where $n \in \mathbb{N}$ is some fixed number. We further assume that there exists a functional $J : \mathcal{D} \rightarrow \mathbb{R}$, such that the desired optical property is characterized by a global minimum of J . Such a functional is referred to as a *goal functional*, henceforth. In the following, we assume that every goal function J is of the form

$$J(\rho) = \Upsilon(\omega_1(\rho, \cdot), \dots, \omega_n(\rho, \cdot)) \quad \text{for all } \rho \in \mathcal{D}, \quad (5.3)$$

where

$$\Upsilon : \underbrace{C^0(\mathbb{B}, \mathbb{R}) \times \dots \times C^0(\mathbb{B}, \mathbb{R})}_{n\text{-times}} \rightarrow \mathbb{R}$$

is a locally Lipschitz continuous functional, which evaluates the first n bands of a band structure with respect to the optimization goal. Recall that the notion of local Lipschitz continuity was introduced in Section 2.3. Notice that by Proposition 4.24 on page 79, we have indeed that the functions $\omega_1(\rho, \cdot), \dots, \omega_n(\rho, \cdot)$ are continuous on \mathbb{B} for every $\rho \in \mathcal{D}$.

With the above definitions, we arrive at the following optimization problem.

Problem 5.1. *Given a goal functional $J : \mathcal{D} \rightarrow \mathbb{R}$ of the form (5.3), find a density function $\rho_* \in \mathcal{C}$, such that*

$$J(\rho_*) \leq J(\rho) \quad \text{for all } \rho \in \mathcal{C}.$$

Clearly, Problem 5.1 constitutes a minimization problem of the form

$$\underset{\rho \in \mathcal{C}}{\text{minimize}} \quad J(\rho).$$

It should be noted that Problem 5.1 is, in fact, only a problem scheme, since the functional Υ was not defined so far. In the following, it is assumed that the term *photonic band structure optimization problem (PBSOP)* refers to an optimization problem, which fits into the scheme of Problem 5.1.

5.2 Regularity of the Goal Functionals

In this section we investigate the regularity of the goal functionals J in Problem 5.1. Recall that these goal functionals are assumed to be of the form given in (5.3). Due to this, we are able to show that the goal functionals are locally Lipschitz continuous on \mathcal{D} and, under further assumptions, even Lipschitz continuous.

The following proposition establishes the local Lipschitz continuity of the eigenvalues $\lambda_j(\rho, \mathbf{k})$ of the weakly formulated eigenvalue equation (4.42) on page 64 with respect to ρ .

Theorem 5.2. *For every index $j \in \mathbb{N}$ and every vector $\mathbf{k} \in \mathbb{B}$ the functional $\lambda_j(\cdot, \mathbf{k}) : \mathcal{D} \rightarrow \mathbb{R}$ is locally Lipschitz continuous.*

Proof. In the following we assume that $j \in \mathbb{N}$ and $\mathbf{k} \in \mathbb{B}$ are fixed. One easily verifies that the mapping

$$\rho \mapsto \operatorname{ess\,inf}_{\Omega}(\rho)$$

is Lipschitz continuous as a functional on \mathcal{E} with Lipschitz constant 1. This implies that the functions $\alpha(\cdot, \mathbf{k})$ and $\kappa(\cdot, \mathbf{k})$ are Lipschitz continuous on \mathcal{D} . Recall that the functions α and κ were defined in (4.38) and (4.40) on page 63. Next, we chose an arbitrary coefficient $\rho_0 \in \mathcal{D}$. Letting

$$\delta_1 := \frac{1}{2} \operatorname{ess\,inf}_{\Omega}(\rho_0),$$

we denote by $\mathcal{B}_{\delta_1}(\rho_0)$ the open ball in \mathcal{E} with radius δ_1 and center ρ_0 . By definition of the function α and the estimate given in (4.23) on page 51, we have that

$$\begin{aligned} \alpha(\rho, \mathbf{k}) &\geq \alpha(\rho_0, \mathbf{k}) - \|\rho - \rho_0\|_{\Omega, \infty} \left(1 - \frac{|\mathbf{k}|}{\operatorname{diam}(\mathbb{B})}\right) \\ &\geq \frac{\delta_1}{4} > 0 \quad \text{for all } \rho \in \mathcal{B}_{\delta_1}(\rho_0). \end{aligned}$$

Given two arbitrary coefficients $\rho_1, \rho_2 \in \mathcal{B}_{\delta_1}(\rho_0)$, we now consider the operators $G_{\mathbf{k}}(\rho_1)$ and $G_{\mathbf{k}}(\rho_2)$ defined by (4.51) on page 66. Let $\mathbf{f} \in \mathbf{Z}$ be an arbitrarily chosen function, and let $\mathbf{u} := L_{\mathbf{k}}(\rho_1)^{-1}M_{\mathbf{k}}\mathbf{f}$ and $\mathbf{v} := L_{\mathbf{k}}(\rho_2)^{-1}M_{\mathbf{k}}\mathbf{f}$. Recall that the operators $M_{\mathbf{k}}$ and $L_{\mathbf{k}}(\rho_1)$ were defined in (4.44) and (4.48) on page 65. It follows that \mathbf{u} and \mathbf{v} are vector fields in $\mathbf{V}_{\mathbf{k}}$. By the norm estimates (4.46) and (4.50) on page 65 and the above estimate for the function α , we have that \mathbf{u} and \mathbf{v} satisfy $\|\mathbf{u}\|_{\operatorname{curl}, \Omega} \|\mathbf{v}\|_{\operatorname{curl}, \Omega} \leq (16/\delta_1^2) \|\mathbf{f}\|_{\Omega}^2$. Since $G_{\mathbf{k}}(\rho_1)$ is self-adjoint, we also have that $\langle G_{\mathbf{k}}(\rho_1)\mathbf{f}, \mathbf{f} \rangle_{\Omega}$ is a real number. Therefore, we obtain the identity

$$\begin{aligned} \langle G_{\mathbf{k}}(\rho_1)\mathbf{f}, \mathbf{f} \rangle_{\Omega} &= \langle I_{\mathbf{k}}L_{\mathbf{k}}(\rho_1)^{-1}M_{\mathbf{k}}\mathbf{f}, \mathbf{f} \rangle_{\Omega} = \langle L_{\mathbf{k}}(\rho_1)^{-1}M_{\mathbf{k}}\mathbf{f}, \mathbf{f} \rangle_{\Omega} \\ &= \langle \mathbf{u}, \mathbf{f} \rangle_{\Omega} = M_{\mathbf{k}}\mathbf{u}[\mathbf{f}] = M_{\mathbf{k}}\mathbf{f}[\mathbf{u}] \\ &= L_{\mathbf{k}}(\rho_2)\mathbf{v}[\mathbf{u}] = L_{\mathbf{k}}(\rho_2)\mathbf{u}[\mathbf{v}] \\ &= a_{\mathbf{k}}(\rho_2)(\mathbf{u}, \mathbf{v}) + \kappa(\rho_2, \mathbf{k})m(\mathbf{u}, \mathbf{v}). \end{aligned}$$

Similarly, we obtain that $\langle G_{\mathbf{k}}(\rho_2)\mathbf{f}, \mathbf{f} \rangle_{\Omega} = a_{\mathbf{k}}(\rho_1)(\mathbf{u}, \mathbf{v}) + \kappa(\rho_1, \mathbf{k})m(\mathbf{u}, \mathbf{v})$, and from this we deduce the estimate

$$\begin{aligned} \left| \langle (G_{\mathbf{k}}(\rho_2) - G_{\mathbf{k}}(\rho_1))\mathbf{f}, \mathbf{f} \rangle_{\Omega} \right| &\leq \left| a_{\mathbf{k}}(\rho_1)(\mathbf{u}, \mathbf{v}) - a_{\mathbf{k}}(\rho_2)(\mathbf{u}, \mathbf{v}) \right| \\ &\quad + \left| \kappa(\rho_1, \mathbf{k}) - \kappa(\rho_2, \mathbf{k}) \right| \left| m(\mathbf{u}, \mathbf{v}) \right|. \end{aligned}$$

Furthermore, by Lemma 4.3 on page 50, we have that

$$\begin{aligned} |a_{\mathbf{k}}(\rho_1)(\mathbf{u}, \mathbf{v}) - a_{\mathbf{k}}(\rho_2)(\mathbf{u}, \mathbf{v})| &= |(a_{\mathbf{k}}(\rho_1 - \rho_2)(\mathbf{u}, \mathbf{v}))| \\ &\leq \beta_0(\mathbf{k}) \|\rho_1 - \rho_2\|_{\Omega, \infty} \|\mathbf{u}\|_{\text{curl}, \Omega} \|\mathbf{v}\|_{\text{curl}, \Omega}. \end{aligned}$$

Since $\kappa(\cdot, \mathbf{k})$ is Lipschitz continuous on \mathcal{D} , there exists a constant $c_1 > 0$, such that $|\kappa(\rho_1, \mathbf{k}) - \kappa(\rho_2, \mathbf{k})| \leq c_1 \|\rho_1 - \rho_2\|_{\Omega, \infty}$. Moreover, we have that $|m(\mathbf{u}, \mathbf{v})| \leq \|\mathbf{u}\|_{\Omega} \|\mathbf{v}\|_{\Omega} \leq \|\mathbf{u}\|_{\text{curl}, \Omega} \|\mathbf{v}\|_{\text{curl}, \Omega}$, so that we finally obtain

$$|\langle (G_{\mathbf{k}}(\rho_2) - G_{\mathbf{k}}(\rho_1)) \mathbf{f}, \mathbf{f} \rangle_{\Omega}| \leq c_2 \|\rho_1 - \rho_2\|_{\Omega, \infty} \|\mathbf{f}\|_{\Omega}^2,$$

where $c_2 := 16(\beta_0(\mathbf{k}) + c_1)/\delta_1^2$. From this inequality we deduce that

$$\begin{aligned} \frac{\langle G_{\mathbf{k}}(\rho_1) \mathbf{f}, \mathbf{f} \rangle_{\Omega}}{\|\mathbf{f}\|_{\Omega}^2} &\geq \frac{\langle G_{\mathbf{k}}(\rho_2) \mathbf{f}, \mathbf{f} \rangle_{\Omega}}{\|\mathbf{f}\|_{\Omega}^2} - c_2 \|\rho_1 - \rho_2\|_{\Omega, \infty} && \text{for all } \mathbf{f} \in \mathbf{Z} \setminus \{\mathbf{0}\}, \\ \frac{\langle G_{\mathbf{k}}(\rho_1) \mathbf{f}, \mathbf{f} \rangle_{\Omega}}{\|\mathbf{f}\|_{\Omega}^2} &\leq \frac{\langle G_{\mathbf{k}}(\rho_2) \mathbf{f}, \mathbf{f} \rangle_{\Omega}}{\|\mathbf{f}\|_{\Omega}^2} + c_2 \|\rho_1 - \rho_2\|_{\Omega, \infty} && \text{for all } \mathbf{f} \in \mathbf{Z} \setminus \{\mathbf{0}\}. \end{aligned}$$

By applying the Min-Max principle, stated in (4.54) on page 69, we find that the j -th largest positive eigenvalues of $G_{\mathbf{k}}(\rho_1)$ and $G_{\mathbf{k}}(\rho_2)$ satisfy $\mu_j(\rho_1, \mathbf{k}) \geq \mu_j(\rho_2, \mathbf{k}) - c_2 \|\rho_1 - \rho_2\|_{\Omega, \infty}$ and $\mu_j(\rho_1, \mathbf{k}) \leq \mu_j(\rho_2, \mathbf{k}) + c_2 \|\rho_1 - \rho_2\|_{\Omega, \infty}$ for every $j \in \mathbb{N}$. Combining both estimates we obtain the inequality $|\mu_j(\rho_1, \mathbf{k}) - \mu_j(\rho_2, \mathbf{k})| \leq c_2 \|\rho_1 - \rho_2\|_{\Omega, \infty}$. Since this inequality can be established for all coefficients $\rho_1, \rho_2 \in \mathcal{B}_{\delta_1}(\rho_0)$, we have that the function $\mu_j(\cdot, \mathbf{k})$ is Lipschitz continuous on $\mathcal{B}_{\delta_1}(\rho_0)$. Setting

$$\delta_2 := \min \left\{ \delta_1, \frac{\mu_j(\rho_0, \mathbf{k})}{2c_2} \right\}$$

we obtain that

$$\mu_j(\rho_1, \mathbf{k}) \geq \mu_j(\rho_0, \mathbf{k}) - c_2 \|\rho_1 - \rho_0\|_{\Omega, \infty} \geq \frac{\mu_j(\rho_0, \mathbf{k})}{2} > 0 \quad \text{for all } \rho_1 \in \mathcal{B}_{\delta_2}(\rho_0).$$

Hence, the function $\mu_j(\cdot, \mathbf{k})$ is uniformly bounded away from zero by $\mu_j(\rho_0, \mathbf{k})/2$ on $\mathcal{B}_{\delta_2}(\rho_0)$. Using the Mean Value Theorem, one easily verifies that the function $1/\mu_j(\cdot, \mathbf{k})$ is Lipschitz continuous on $\mathcal{B}_{\delta_2}(\rho_0)$. By (4.60) on page 70 the function $\lambda_j(\cdot, \mathbf{k})$ is given as the sum of the functions $1/\mu_j(\cdot, \mathbf{k})$ and $-\kappa(\cdot, \mathbf{k})$. Hence, $\lambda_j(\cdot, \mathbf{k})$ is Lipschitz continuous on $\mathcal{B}_{\delta_2}(\rho_0)$, which establishes the assertion of this theorem. \square

We can prove an even stronger regularity result for restrictions of the functions $\lambda_j(\cdot, \mathbf{k})$ to a specific subset of \mathcal{D} . Before, however, we state the following, auxiliary result.

Proposition 5.3. *For every index $j \in \mathbb{N}$ and every vector $\mathbf{k} \in \mathbb{B}$ the functional $\lambda_j(\cdot, \mathbf{k}) : \mathcal{D} \rightarrow \mathbb{R}$ is monotonically increasing, i.e., for every two coefficients $\rho_1, \rho_2 \in \mathcal{D}$ that satisfy $\rho_1 \leq \rho_2$ almost everywhere on Ω , we have*

$$\lambda_j(\rho_1, \mathbf{k}) \leq \lambda_j(\rho_2, \mathbf{k}).$$

Proof. Let $\rho_1, \rho_2 \in \mathcal{D}$, such that $\rho_1 \leq \rho_2$ almost everywhere on Ω , and let $\mathbf{u} \in \mathbf{V}_{\mathbf{k}}$. Then, it follows immediately from the definition of the sesquilinear forms $a_{\mathbf{k}}(\rho)$, which was given in (4.33) on page 59, that $a_{\mathbf{k}}(\rho_1)(\mathbf{u}, \mathbf{u}) \leq a_{\mathbf{k}}(\rho_2)(\mathbf{u}, \mathbf{u})$. Since this inequality can be established for arbitrary functions $\mathbf{u} \in \mathbf{V}_{\mathbf{k}}$, we have that

$$\frac{a_{\mathbf{k}}(\rho_1)(\mathbf{u}, \mathbf{u})}{m(\mathbf{u}, \mathbf{u})} \leq \frac{a_{\mathbf{k}}(\rho_2)(\mathbf{u}, \mathbf{u})}{m(\mathbf{u}, \mathbf{u})} \quad \text{for all } \mathbf{u} \in \mathbf{V}_{\mathbf{k}} \setminus \{\mathbf{0}\}.$$

Using the Min-Max Principle stated in Proposition 4.20 on page 70, the above inequality implies that $\lambda_j(\rho_1, \mathbf{k}) \leq \lambda_j(\rho_2, \mathbf{k})$. \square

So far we considered the j -th smallest eigenvalue $\lambda_j(\cdot, \mathbf{k})$ as a functional on the set \mathcal{D} for given $j \in \mathbb{N}$ and given $\mathbf{k} \in \mathbb{N}$. Next, we consider the restrictions of this functional to the function set

$$\bar{\mathcal{C}} := \{\rho \in \mathcal{E} \mid \rho_{\min} \leq \rho \leq \rho_{\max} \text{ almost everywhere in } \Omega\}, \quad (5.4)$$

where the constants ρ_{\min} and ρ_{\max} coincide with those defining the admissible set \mathcal{C} according to (5.2). We shall refer to the set $\bar{\mathcal{C}}$ as the *extended admissible set*. Functions, which belong to $\bar{\mathcal{C}}$, are also referred to as density functions. One easily verifies that $\mathcal{C} \subset \bar{\mathcal{C}} \subset \mathcal{D}$. With the extended admissible set $\bar{\mathcal{C}}$ at hand, we can state the following theorem.

Theorem 5.4. *For every index $j \in \mathbb{N}$ and every vector $\mathbf{k} \in \mathbb{B}$ the functional $\lambda_j(\cdot, \mathbf{k})$ is Lipschitz continuous on $\bar{\mathcal{C}}$.*

Proof. Clearly, every density function $\rho \in \rho_{\min}$ is positive and uniformly bounded away from zero by ρ_{\min} almost everywhere in Ω . Therefore, we have that $\alpha(\rho, \mathbf{k}) > \rho_{\min}/2$ for all $\rho \in \bar{\mathcal{C}}$. By the same line of arguments as in the proof of Theorem 5.2, we deduce that the function $\mu_j(\cdot, \mathbf{k})$ is Lipschitz continuous on $\bar{\mathcal{C}}$. By Proposition 5.3 the corresponding function $\lambda_j(\cdot, \mathbf{k})$ is monotonically increasing on \mathcal{D} . One easily verifies, that the same holds true for the function $\kappa(\cdot, \mathbf{k})$. Since $1/\mu_j(\cdot, \mathbf{k})$ coincides with $\lambda_j(\cdot, \mathbf{k}) + \kappa(\cdot, \mathbf{k})$ by (4.60) on page 70, we conclude that $\mu_j(\cdot, \mathbf{k})$ is monotonically decreasing on $\bar{\mathcal{C}}$. Hence, $\mu_j(\cdot, \mathbf{k})$ is uniformly bounded from below by $\mu_j(\rho_{\min}, \mathbf{k}) > 0$ on $\bar{\mathcal{C}}$. Because of this, the functional $1/\mu_j(\cdot, \mathbf{k})$, and hence also the functional $\lambda_j(\cdot, \mathbf{k})$, are Lipschitz continuous on $\bar{\mathcal{C}}$. \square

We now relate the result of Theorem 5.2 to the regularity of the goal functional J in Problem 5.1. According to (4.74) on page 80, the functionals ω_j are given by $\omega_j := \sqrt{\lambda_j}$ for every index $j \in \mathbb{N}$. Using the Mean-Value Theorem, one can show that the function $x \mapsto \sqrt{x}$ is locally Lipschitz continuous on $\mathbb{R}_{>0}$. Since the composition of two locally Lipschitz continuous functions is locally Lipschitz continuous again (see Section 2.3), we deduce that the functions $\omega_j(\cdot, \mathbf{k})$ are locally Lipschitz continuous for every index $j \in \mathbb{N}$ and every vector $\mathbf{k} \in \mathbb{B}$. Since the functional Υ , defined in (5.3), is also locally Lipschitz continuous by assumption, we deduce the following corollary to Theorem 5.2.

Corollary 5.5. *The goal functional $J : \mathcal{D} \rightarrow \mathbb{R}$ in Problem 5.1 is locally Lipschitz continuous.*

Finally, we remark that under further assumptions on the functional Υ , the goal functional J can be shown to be Lipschitz continuous on the set $\bar{\mathcal{C}}$. Consider, for example, the case where the optimization problem only involves the first photonic band, and where the functional $\Upsilon : C^0(\mathbb{B}, \mathbb{R}) \rightarrow \mathbb{R}$ is given by

$$\Upsilon(f) = f(\mathbf{k}_0) \quad \text{for all } f \in C^0(\mathbb{B}, \mathbb{R}),$$

for some prescribed vector $\mathbf{k}_0 \in \mathbb{B} \setminus \{\mathbf{0}\}$. According to Proposition 4.20 on page 70 and Proposition 5.3 we that

$$\lambda_1(\rho, \mathbf{k}_0) \geq \lambda_1(\rho_{\min}, \mathbf{k}_0) > 0 \quad \text{for all } \rho \in \bar{\mathcal{C}}.$$

Hence, the function $\lambda_1(\cdot, \mathbf{k}_0)$ is uniformly bounded away from zero on the extended admissible set $\bar{\mathcal{C}}$ by $\lambda_1(\rho_{\min}, \mathbf{k}_0)$. Since $\lambda_1(\cdot, \mathbf{k}_0)$ is Lipschitz continuous on $\bar{\mathcal{C}}$ according to Theorem 5.4, it follows that the function $\omega_1(\cdot, \mathbf{k}_0) = \sqrt{\lambda_1(\cdot, \mathbf{k}_0)}$, and hence also the goal functional $J = \omega_1(\cdot, \mathbf{k}_0)$, are Lipschitz continuous on $\bar{\mathcal{C}}$.

5.3 Existence of Optima in the TM Setting

So far, we considered photonic band structure optimization problems (PBSOPs) for three-dimensional photonic crystals. In this section we shift the focus of our investigations to PBSOPs, which involve TM band structures of two-dimensional photonic crystals.

Throughout this section, we assume that $\Lambda \subset \mathbb{R}^2$ is a Bravais lattice of rank 2, and that Ω is a primitive domain of Λ . By \mathbb{B} we denote the first Brillouin zone of Λ . In Section 4.7 we defined $\lambda_j^{\text{TM}}(\rho, \mathbf{k})$ as the j -th smallest eigenvalue of the eigenvalue equation (4.83) on page 85 for given $\rho \in \mathcal{D}$ and given $\mathbf{k} \in \mathbb{B}$. The sequence of the corresponding eigenfunctions $\{u_j^{\text{TM}}(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$ forms a complete system of V . Recall

that V denotes the complex Sobolev space $H_{\text{per}}^1(\Omega)$. The system is orthogonal with respect to the inner product $\langle \cdot, \cdot \rangle_\rho$, which was defined by (4.87) on page 86.

In analogy to the three-dimensional setting defined in the Sections 5.1 and Sections 5.2, we define the function set \mathcal{C} as

$$\mathcal{C} := \{ \rho_{\min}(1 - \chi_S) + \rho_{\max}\chi_S \mid S \in \mathcal{M}(\Omega) \},$$

and the function set $\bar{\mathcal{C}}$ as

$$\bar{\mathcal{C}} := \{ \rho \in \mathcal{E} \mid \rho_{\min} \leq \rho \leq \rho_{\max} \text{ almost everywhere in } \Omega \}, \quad (5.5)$$

for some constants $\rho_{\min}, \rho_{\max} \in \mathbb{R}$ satisfying $0 < \rho_{\min} < \rho_{\max}$. As in the three-dimensional setting, the sets \mathcal{C} is called the *admissible set*, and the set $\bar{\mathcal{C}}$ is called the *extended admissible set*. Clearly, every function $\rho \in \mathcal{C}$ represents the medium structure of a two-dimensional photonic crystal consisting of exactly two different materials. The Λ -periodic relative electric permittivity function ε_r of such a crystal is determined by $\varepsilon_r|_\Omega = \rho$.

Given some fixed number $n \in \mathbb{N}$, a functional $J^{\text{TM}} : \mathcal{D} \rightarrow \mathbb{R}$, which is given by

$$J^{\text{TM}}(\rho) = \Upsilon^{\text{TM}}(\omega_1^{\text{TM}}(\rho, \cdot), \dots, \omega_n^{\text{TM}}(\rho, \cdot)) \quad \text{for all } \rho \in \mathcal{D}, \quad (5.6)$$

for some locally Lipschitz continuous functional

$$\Upsilon^{\text{TM}} : \underbrace{C^0(\mathbb{B}, \mathbb{R}) \times \dots \times C^0(\mathbb{B}, \mathbb{R})}_{n\text{-times}} \rightarrow \mathbb{R},$$

is called a *TM goal functional*. In analogy to the three-dimensional setting (see Section 5.1), we assume that the TM goal functional J^{TM} attains its minimum at some density function $\rho_* \in \mathcal{C}$, if and only if the TM band structure of the two-dimensional photonic crystal, whose medium structure is represented by ρ_* , exhibits some desired property. We then consider the following problem scheme.

Problem 5.6. *Given a TM goal functional $J^{\text{TM}} : \mathcal{D} \rightarrow \mathbb{R}$ of the form (5.6), find an density function $\rho_* \in \mathcal{C}$, such that*

$$J^{\text{TM}}(\rho_*) \leq J^{\text{TM}}(\rho) \quad \text{for all } \rho \in \mathcal{C}.$$

Clearly, Problem 5.6 is an analogue of Problem 5.1 for two-dimensional photonic crystals, and by a similar line of argument as in Section 5.2, one can show that every TM goal functional is locally Lipschitz continuous on \mathcal{D} .

In the remainder of this section we shall address the question as to whether or not the optimization problem stated in Problem 5.6 admits a solution. We remark

that similar optimization problems arise in the design of frequency-optimized, two-composite membranes. Such problems have been studied extensively by Cox and McLaughling (cf. [31], [32]). It turns out, that one has to study such problems on the extended admissible set $\bar{\mathcal{C}}$ to prove the existence of solutions.

In the following, we introduce a topology on the function space \mathcal{E} that is weaker than the norm topology. Recall that $L^1(\Omega, \mathbb{R})^*$, the normed dual of the real Banach space $L^1(\Omega, \mathbb{R})$, is isometrically linearly isomorphic to $\mathcal{E} = L^\infty(\Omega, \mathbb{R})$ (see e.g. Theorem 4.12 in [4]). The canonical isometric, linear isomorphism $\Phi : L^\infty(\Omega, \mathbb{R}) \rightarrow L^1(\Omega, \mathbb{R})^*$ is given by

$$\Phi(\xi)[f] := \int_{\Omega} \xi f \quad \text{for all } \xi \in \mathcal{E}, f \in L^1(\Omega, \mathbb{R}). \quad (5.7)$$

The *weak*-topology on $L^1(\Omega, \mathbb{R})^*$* is defined as the weakest topology on $L^1(\Omega, \mathbb{R})^*$, with respect to which all continuous linear functionals $\varphi : L^1(\Omega, \mathbb{R})^* \rightarrow \mathbb{R}$, that are given by $\varphi(p) = p(f)$ for some $f \in L^1(\Omega, \mathbb{R})$ and for all $p \in L^1(\Omega, \mathbb{R})^*$, are continuous. By virtue of the isomorphism Φ every open set in that topology can be identified with a subset of \mathcal{E} . Thus, the weak*-topology on $L^1(\Omega, \mathbb{R})^*$ also defines a topology on \mathcal{E} which is called the *weak*-topology on \mathcal{E}* . A sequence in \mathcal{E} , which is convergent with respect to this topology, is called *weak*-convergent* in \mathcal{E} . One easily verifies that a sequence $\{\xi^{(l)}\}_{l \in \mathbb{N}}$ in \mathcal{E} is weak*-convergent in \mathcal{E} , if and only if there exists a function $\xi \in \mathcal{E}$, such that

$$\Phi(\xi^{(l)} - \xi)[f] = \int_{\Omega} (\xi^{(l)} - \xi)f \rightarrow 0 \quad \text{as } l \rightarrow \infty \quad \text{for all } f \in L^1(\Omega, \mathbb{R}).$$

In this case we shall say that the sequence $\{\xi^{(l)}\}_{l \in \mathbb{N}}$ weak*-converges to ξ . Obviously, convergence in \mathcal{E} , i.e., convergence with respect to the norm topology on \mathcal{E} , also implies weak*-convergence in \mathcal{E} . The converse, however, does not hold in general. Hence, the norm topology on \mathcal{E} is strictly stronger than the weak*-topology on \mathcal{E} .

In the following, terms like weak*-compactness or weak*-continuity always refer to the respective properties with respect to the weak*-topology on \mathcal{E} . The following proposition states in particular that the extended admissible set $\bar{\mathcal{C}}$ is compact with respect to that topology.

Proposition 5.7. *The set $\bar{\mathcal{C}}$ is convex, bounded, and weak*-compact.*

Proof. Convexity and boundedness directly follow from the definition of $\bar{\mathcal{C}}$ in (5.5). It hence remains to show the weak*-compactness of $\bar{\mathcal{C}}$. To this end, we denote by $\overline{\mathcal{B}_1(0)}$ the closed unit ball in \mathcal{E} and define the operator $\Psi : \bar{\mathcal{C}} \rightarrow \overline{\mathcal{B}_1(0)}$ by

$$\Psi(\rho) := \frac{2\rho - \rho_{\max} - \rho_{\min}}{\rho_{\max} - \rho_{\min}} \quad \text{for all } \rho \in \bar{\mathcal{C}}.$$

One easily verifies that Ψ is a bijection. Choosing an arbitrary sequence $\{\rho^{(l)}\}_{l \in \mathbb{N}}$ in $\overline{\mathcal{C}}$, we define the sequence $\{\xi^{(l)}\}_{l \in \mathbb{N}}$ in $\overline{\mathcal{B}_1(0)}$ by $\xi^{(l)} := \Psi(\rho^{(l)})$ for all $l \in \mathbb{N}$, as well as the sequence $\{p^{(l)}\}_{l \in \mathbb{N}}$ in $L^1(\Omega, \mathbb{R})^*$ by $p^{(l)} := \Phi(\xi^{(l)})$ for all $l \in \mathbb{N}$. Here, Φ denotes the isometric linear isomorphism given by (5.7). Since

$$\|p^{(l)}\|_{L^1(\Omega, \mathbb{R})^*} = \|\Psi(\rho^{(l)})\|_{\Omega, \infty} \leq 1 \quad \text{for all } l \in \mathbb{N},$$

we have that $\{p^{(l)}\}_{l \in \mathbb{N}}$ is a sequence in the closed unit ball $\overline{B_1^*(0)}$ in $L^1(\Omega, \mathbb{R})^*$. According to Alaoglu's Theorem (see e.g. Theorem 13.9 [40]), $B_1^*(0)$ is compact with respect to the weak-* topology on $L^1(\Omega, \mathbb{R})^*$. Hence, there exists a subsequence $\{p^{(1,l)}\}_{l \in \mathbb{N}}$ of $\{p^{(l)}\}_{l \in \mathbb{N}}$ and a continuous linear functional $p \in L^1(\Omega, \mathbb{R})^*$ such that

$$(p^{(1,l)} - p)(f) \rightarrow 0 \quad \text{as } l \rightarrow \infty \quad \text{for all } f \in L^1(\Omega, \mathbb{R}).$$

Letting $\xi := \Phi^{-1}(p)$ and $\rho := \Psi^{-1}(\xi)$, we have that

$$\int_{\Omega} (\rho^{(1,l)} - \rho)f = \frac{\rho_{\max} - \rho_{\min}}{2} \int_{\Omega} (\xi^{(1,l)} - \xi)f = (p^{(1,l)} - p)(f) \quad \text{for all } l \in \mathbb{N}.$$

Hence, $\{\rho^{(1,l)}\}_{l \in \mathbb{N}}$ is a weak-* convergent subsequence of $\{\rho^{(l)}\}_{l \in \mathbb{N}}$. Since $\{\rho^{(l)}\}_{l \in \mathbb{N}}$ was chosen arbitrarily, it follows that $\overline{\mathcal{C}}$ is compact with respect to the weak-* topology on \mathcal{E} . \square

Next, we study the relationship between the sets \mathcal{C} and $\overline{\mathcal{C}}$ with respect to the weak-* topology on \mathcal{E} . The following theorem is a particularization of Theorem 3 in [72].

Theorem 5.8. *Let R be a non-empty subset of \mathbb{R} .*

- (a) *Let $\{\xi^{(l)}\}_{l \in \mathbb{N}}$ be a sequence in \mathcal{E} , such that $\{\xi^{(l)}\}_{l \in \mathbb{N}}$ weak-* converges to a function $\xi \in \mathcal{E}$, and such that $\xi^{(l)}(\mathbf{x}) \in R$ for almost all $\mathbf{x} \in \Omega$. Then, $\xi(\mathbf{x}) \in \overline{\text{conv}}(R)$ for almost all $\mathbf{x} \in \Omega$.*
- (b) *Conversely, let $\xi \in \mathcal{E}$ be a function, such that $\xi(\mathbf{x}) \in \overline{\text{conv}}(R)$ for almost all $\mathbf{x} \in \Omega$. Then, there exists a sequence $\{\xi^{(l)}\}_{l \in \mathbb{N}}$ in \mathcal{E} , such that $\xi^{(l)}(\mathbf{x}) \in R$ for almost all $\mathbf{x} \in \Omega$ and all $l \in \mathbb{N}$, and such that $\{\xi^{(l)}\}_{l \in \mathbb{N}}$ weak-* converges to ξ .*

Choosing $R = \{\rho_{\min}, \rho_{\max}\}$ in Theorem 5.8, we obtain that every weak-* limit function of a weak-* convergent sequence in the admissible set \mathcal{C} takes function values in the closed interval $[\rho_{\min}, \rho_{\max}]$ almost everywhere on Ω . Therefore, every such limit function belongs to the extended admissible set $\overline{\mathcal{C}}$. Conversely, for every

function in the extended admissible $\bar{\mathcal{C}}$ there exists a sequence in the admissible set \mathcal{C} , which weak- $*$ -converges to that function. Hence, we have the following corollary to Theorem 5.8.

Corollary 5.9. *The set $\bar{\mathcal{C}}$ is the weak- $*$ -closure of \mathcal{C} .*

By definition the TM goal functional J^{TM} is locally Lipschitz continuous and hence continuous on \mathcal{D} with respect to the norm topology on \mathcal{E} . Since the norm topology is strictly stronger than the weak- $*$ -topology on \mathcal{E} , the weak- $*$ -compactness of the extended admissible set $\bar{\mathcal{C}}$ does not imply that the TM goal functional attains a minimum on that set. We first need to establish weak- $*$ -continuity for the TM goal functional, in order to obtain such a result. The following lemma is an auxiliary result and follows from Lemma 4.2 in [31]. Recall that Z denotes the complex Lebesgue space $L^2(\Omega)$.

Lemma 5.10. *Let $\{\rho^{(l)}\}_{l \in \mathbb{N}}$ be a sequence in $\bar{\mathcal{C}}$ that weak- $*$ -converges to $\rho \in \bar{\mathcal{C}}$, and let $\{f^{(l)}\}_{l \in \mathbb{N}}$ and $\{g^{(l)}\}_{l \in \mathbb{N}}$ be sequences in Z that converge in Z to $f \in Z$ and $g \in Z$, respectively. Then, we have that*

$$m^{\text{TM}}(\rho^{(l)})(f^{(l)}, g^{(l)}) \rightarrow m^{\text{TM}}(\rho)(f, g) \quad \text{as } l \rightarrow \infty.$$

The following theorem states that the functions $\lambda_j^{\text{TM}}(\cdot, \mathbf{k})$ are weak- $*$ -continuous.

Theorem 5.11. *Let $\{\rho^{(l)}\}_{l \in \mathbb{N}}$ be a sequence in $\bar{\mathcal{C}}$, which weak- $*$ -converges to a coefficient $\rho \in \bar{\mathcal{C}}$. Then, we have for every index $j \in \mathbb{N}$, and every vector $\mathbf{k} \in \mathbb{B}$ that the sequence $\{\lambda_j^{\text{TM}}(\rho^{(l)}, \mathbf{k})\}_{l \in \mathbb{N}}$ converges to $\lambda_j^{\text{TM}}(\rho, \mathbf{k})$.*

Proof. This proof is an adaptation of the proof of Proposition 4.3(i) in [31]. Let $\mathbf{k} \in \mathbb{B}$ be fixed, and let $\{\rho^{(l)}\}_{l \in \mathbb{N}}$ be a sequence in $\bar{\mathcal{C}}$ that weak- $*$ -converges to a density function $\rho \in \bar{\mathcal{C}}$. In order to simplify notations, we define $\lambda_j^{(l)} := \lambda_j(\rho^{(l)}, \mathbf{k})$ and $u_j^{(l)} := u_j^{\text{TM}}(\rho^{(l)}, \mathbf{k})$, as well as $\lambda_j := \lambda_j^{\text{TM}}(\rho, \mathbf{k})$ and $u_j := u_j^{\text{TM}}(\rho, \mathbf{k})$ for all $j, l \in \mathbb{N}$. Without loss of generality we assume that the sequences $\{u_j^{(l)}\}_{j \in \mathbb{N}}$ and $\{u_j\}_{j \in \mathbb{N}}$ form a complete system of V , which is orthonormal with respect to the inner product $\langle \cdot, \cdot \rangle_\rho$.

Using the Min-Max Principle (4.88) on page 86, one can show that the mapping $\lambda_j^{\text{TM}}(\cdot, \mathbf{k})$ is monotonically decreasing for all $j \in \mathbb{N}$ and all $\mathbf{k} \in \mathbb{B}$. Hence, for every index $j \in \mathbb{N}$, the sequence $\{\lambda_j^{(l)}\}_{l \in \mathbb{N}}$ is bounded from above by $\lambda_j^{\text{TM}}(\rho_{\min}, \mathbf{k})$ and from below by $\lambda_j^{\text{TM}}(\rho_{\max}, \mathbf{k})$. Since the sesquilinear forms $a_{\mathbf{k}}^{\text{TM}}$ are coercive, we have that

$$\alpha^{\text{TM}}(\mathbf{k}) \|u_j^{(l)}\|_{1, \Omega}^2 \leq a_{\mathbf{k}}^{\text{TM}}(u_j^{(l)}, u_j^{(l)}) + \kappa^{\text{TM}}(\mathbf{k}) \|u_j^{(l)}\|_{\Omega}^2$$

$$\begin{aligned}
&\leq \lambda_j^{\text{TM}}(\rho^{(l)}, \mathbf{k}) + \kappa^{\text{TM}}(\mathbf{k}) \int_{\Omega} \frac{\rho^{(l)}}{\rho_{\min}} |u_j^{(l)}|^2 \\
&\leq \lambda_j^{\text{TM}}(\rho_{\min}, \mathbf{k}) + \frac{\kappa^{\text{TM}}(\mathbf{k}) \|u_j^{(l)}\|_{\rho}^2}{\rho_{\min}} \\
&= \lambda_j^{\text{TM}}(\rho_{\min}, \mathbf{k}) + \frac{\kappa^{\text{TM}}(\mathbf{k}) \|u_j^{(l)}\|_{\rho}^2}{\rho_{\min}} \quad \text{for all } j, l \in \mathbb{N}.
\end{aligned}$$

Recall that the functions α^{TM} and κ^{TM} were defined by (4.84) and (4.84) on page 85, respectively. From the above estimate we deduce that, for every index $j \in \mathbb{N}$, the sequence $\{u_j^{(l)}\}_{l \in \mathbb{N}}$ is bounded in V .

It follows that there exists a subsequence $\{\rho^{(1,l)}\}_{l \in \mathbb{N}}$ of $\{\rho^{(l)}\}_{l \in \mathbb{N}}$, such that the corresponding subsequence $\{\lambda_1^{(1,l)}\}_{l \in \mathbb{N}}$ of $\{\lambda_1^{(l)}\}_{l \in \mathbb{N}}$ converges to a non-negative, real number $\nu_1 \in \mathbb{R}_{\geq 0}$, and such that the subsequence $\{u_1^{(1,l)}\}_{l \in \mathbb{N}}$ of $\{u_1^{(l)}\}_{l \in \mathbb{N}}$ converges weakly in V to a function $w_1 \in V$. The existence of the first subsequence follows from the Theorem of Bolzano–Weierstrass. The existence of the second subsequence follows from the fact that every bounded subset of a Hilbert space is weakly sequentially compact (see e.g. Theorem 6.9 in [4]). Since the identical embedding from $H^1(\Omega)$ into Z is compact, we also have that the subsequence $\{u_1^{(1,l)}\}_{l \in \mathbb{N}}$ converges strongly in Z to w_1 (see e.g. Lemma 8.2 in [4]). By the same reasoning, we can extract a second subsequence $\{\rho^{(2,l)}\}_{l \in \mathbb{N}}$ of $\{\rho^{(1,l)}\}_{l \in \mathbb{N}}$, such that $\{\lambda_2^{(2,l)}\}_{l \in \mathbb{N}}$ converges to a non-negative, real number $\nu_2 \in \mathbb{R}_{\geq 0}$, and such that $\{u_2^{(2,l)}\}_{l \in \mathbb{N}}$ converges weakly in V and strongly in Z to a function $w_2 \in V$. By repeating this procedure, we finally obtain a sequence $\{\{\rho^{(j,l)}\}_{l \in \mathbb{N}}\}_{j \in \mathbb{N}}$ of sequences in \mathcal{C} with the property that $\{\rho^{(j+1,l)}\}_{l \in \mathbb{N}}$ is a subsequence of $\{\rho^{(j,l)}\}_{l \in \mathbb{N}}$ for all $j \in \mathbb{N}$. Furthermore, we have for every index $j \in \mathbb{N}$ that the sequence $\{\lambda_j^{(j,l)}\}_{l \in \mathbb{N}}$ converges to a non-negative, real number $\nu_j \in \mathbb{R}_{\geq 0}$, and that the sequence $\{u_j^{(j,l)}\}_{l \in \mathbb{N}}$ converges weakly in V and strongly in Z to a function $w_j \in V$. We now consider the “diagonal” subsequence $\{\rho^{(l,l)}\}_{l \in \mathbb{N}}$ of $\{\rho^{(l)}\}_{l \in \mathbb{N}}$. It follows from the discussion above, that for every index $j \in \mathbb{N}$ the sequence $\{\lambda_j^{(l,l)}\}_{l \in \mathbb{N}}$ converges to ν_j . Furthermore, we have that the sequence $\{u_j^{(l,l)}\}_{l \in \mathbb{N}}$ converges weakly in V and strongly in Z to w_j .

By Lemma 5.10, we have that the sequence $\{w_j\}_{j \in \mathbb{N}}$ forms an orthonormal system with respect to the inner product $\langle \cdot, \cdot \rangle_{\rho}$, since

$$\begin{aligned}
\delta_{ij} &= \langle u_i^{(l,l)}, u_j^{(l,l)} \rangle_{\rho^{(l)}} = m^{\text{TM}}(\rho^{(l,l)})(u_i^{(l,l)}, u_j^{(l,l)}) \\
&\rightarrow m^{\text{TM}}(\rho)(w_i, w_j) = \langle w_i, w_j \rangle_{\rho} \quad \text{as } l \rightarrow \infty \quad \text{for all } i, j \in \mathbb{N}.
\end{aligned}$$

Furthermore, since $\{\lambda_j^{(l,l)}\}_{j \in \mathbb{N}}$ is an increasing sequence for all $l \in \mathbb{N}$, we also have that $\{\nu_j\}_{j \in \mathbb{N}}$ is an increasing sequence. Since every sequence $\{\lambda_j^{(l,l)}\}_{j \in \mathbb{N}}$ is unbounded, we also have that $\nu_j \rightarrow \infty$ as $j \rightarrow \infty$.

Our next aim is to show that for every index $j \in \mathbb{N}$ there exists an index $i \in \mathbb{N}$, such that $\nu_j = \lambda_i$. We choose $j \in \mathbb{N}$ arbitrarily. Due to the weak converge of the sequence $\{u_j^{(l,l)}\}_{l \in \mathbb{N}}$ in V , we have that

$$a_{\mathbf{k}}^{\text{TM}}(u_j^{(l,l)}, v) \rightarrow a_{\mathbf{k}}^{\text{TM}}(w_j, v) \quad \text{as } l \rightarrow \infty \quad \text{for all } v \in V.$$

By Lemma 5.10, we also have that

$$m^{\text{TM}}(\rho^{(l,l)})(u_j^{(l,l)}, v) \rightarrow m^{\text{TM}}(\rho)(w_j, v) \quad \text{as } l \rightarrow \infty. \quad \text{for all } v \in V.$$

Since $\lambda_j^{(l,l)} \rightarrow \nu_j$ as $l \rightarrow \infty$, we thus obtain

$$a_{\mathbf{k}}^{\text{TM}}(\rho)(w_j, v) = \nu_j m(w_j, v) \quad \text{for all } v \in V,$$

It follows that ν_j is an eigenvalue, and that w_j is a corresponding eigenfunction. Hence, there exists an index $i \in \mathbb{N}$, such that $\nu_j = \lambda_i$.

Next, we show that for every index $i \in \mathbb{N}$ there exists an index $j \in \mathbb{N}$, such that $\lambda_i = \nu_j$. We proof this by contradiction. Suppose that there exists an index $i \in \mathbb{N}$, such that $\lambda_i \neq \nu_j$ for all $j \in \mathbb{N}$. Then, we define

$$\tilde{\lambda} := \lambda_i, \quad \tilde{u} := \frac{u_i}{\tilde{\lambda} + \kappa^{\text{TM}}(\mathbf{k})}.$$

Given an arbitrary index $j \in \mathbb{N}$, we find that

$$\nu_j m^{\text{TM}}(\rho)(w_j, \tilde{u}) = a_{\mathbf{k}}^{\text{TM}}(w_j, \tilde{u}) = \overline{a_{\mathbf{k}}^{\text{TM}}(\tilde{u}, w_j)} = \tilde{\lambda} \overline{m^{\text{TM}}(\rho)(\tilde{u}, w_j)}.$$

Since $\tilde{\lambda} \neq \nu_j$ for all $j \in \mathbb{N}$ by assumption, we conclude that $m^{\text{TM}}(\rho)(\tilde{u}, w_j) = \langle \tilde{u}, w_j \rangle_{\rho} = 0$ for all $j \in \mathbb{N}$. Because the sequence $\{u_j^{(l,l)}\}_{l \in \mathbb{N}}$ converges strongly in Z to w_j for every index $j \in \mathbb{N}$, we have by Lemma 5.10 that

$$P_{\mathbf{k}, j-1}^{\text{TM}}(\rho^{(l,l)})(\tilde{u}) = \sum_{i=1}^{j-1} \langle \tilde{u}, u_i^{(l,l)} \rangle_{\rho^{(l,l)}} u_i^{(l,l)} \rightarrow \sum_{i=1}^{j-1} \langle \tilde{u}, w_i \rangle_{\rho} w_i = 0 \quad \text{as } l \rightarrow \infty.$$

By Auchmuty's principle, stated in (4.91) on page 86, we have that

$$\frac{-1}{2(\lambda_j^{(l,l)} - \kappa^{\text{TM}}(\mathbf{k}))} \leq \left(\frac{1}{2} (a_{\mathbf{k}}^{\text{TM}}(\tilde{u}, \tilde{u}) + \kappa^{\text{TM}}(\mathbf{k}) m^{\text{TM}}(\rho^{(l,l)})(\tilde{u}, \tilde{u})) - \|\tilde{u} - P_{\mathbf{k}, j-1}^{\text{TM}}(\rho^{(l,l)})\tilde{u}\|_{\rho^{(l,l)}} \right) \quad \text{for all } j, l \in \mathbb{N}.$$

Letting l tend to infinity in the above inequality, we obtain

$$\frac{-1}{2(\nu_j - \kappa^{\text{TM}}(\mathbf{k}))} \leq \left(\frac{1}{2} (a_{\mathbf{k}}^{\text{TM}}(\tilde{u}, \tilde{u}) + \kappa^{\text{TM}}(\mathbf{k}) m^{\text{TM}}(\rho)(\tilde{u}, \tilde{u})) \right)$$

$$\begin{aligned}
& - \|\tilde{u} - P_{\mathbf{k},j-1}^{\text{TM}}(\rho)\tilde{u}\|_{\rho} \Big) \\
& = \left(\frac{1}{2} (\tilde{\lambda} + \kappa^{\text{TM}}(\mathbf{k})) \|\tilde{u}\|_{\rho}^2 - \|\tilde{u}\|_{\rho} \right) \\
& = \frac{-1}{2(\tilde{\lambda} - \kappa^{\text{TM}}(\mathbf{k}))} \quad \text{for all } j \in \mathbb{N}.
\end{aligned}$$

Notice that this implies $\nu_j \leq \tilde{\lambda}$ for all $j \in \mathbb{N}$, which is a contradiction to the fact that $\nu_j \rightarrow \infty$ as $j \rightarrow \infty$. We thus have that for every $i \in \mathbb{N}$ there exists a $j \in \mathbb{N}$, such that $\nu_i = \lambda_j$. It follows that the sequences $\{\nu_j\}_{j \in \mathbb{N}}$ and $\{\lambda_j^{\text{TM}}(\rho, \mathbf{k})\}_{j \in \mathbb{N}}$ coincide. \square

Clearly, the weak-* continuity of the functions $\lambda_j^{\text{TM}}(\cdot, \mathbf{k})$ also implies the weak-* continuity of the functions $\omega_j^{\text{TM}}(\cdot, \mathbf{k})$ on $\bar{\mathcal{C}}$, where $j \in \mathbb{N}$. It thus follows from (5.6) that the TM goal functional J^{TM} is also weak-* continuous on $\bar{\mathcal{C}}$. Since $\bar{\mathcal{C}}$ is weak-* compact according to Proposition 5.7, we obtain the following corollary to Theorem 5.11.

Corollary 5.12. *The TM goal functional J^{TM} in Problem 5.6 attains a minimum on the extended admissible set $\bar{\mathcal{C}}$.*

It should be noted that Corollary 5.12 only guarantees the existence of optimal solutions in the extended admissible set $\bar{\mathcal{C}}$. In view of PBSOPs problems, it would be desirable to have optimal solutions in the admissible set \mathcal{C} . This is because only solutions which belong to the admissible set represent two-valued relative electric permittivity functions. According to Proposition 5.7, the extended admissible set $\bar{\mathcal{C}}$ is a convex set. It can be shown that the set of extremal points of $\bar{\mathcal{C}}$ is precisely the admissible set \mathcal{C} (see Proposition 2.5 in [31]). Therefore, every convex and weak-* continuous functional on $\bar{\mathcal{C}}$ attains its minimum on \mathcal{C} . In Section 5.5 we introduce a special class of TM goal functionals. Unfortunately, we could not show that these functionals are convex. The numerical results we present in Chapter 9 indicate, however, that these functionals indeed attain their minima on the admissible set.

5.4 Existence of Optima in Other Settings

In the previous section, we showed that PBSOPs involving the TM band structure of a two-dimensional photonic crystal admit an optimal solution in the extended admissible set $\bar{\mathcal{C}}$. This set consists of all functions in \mathcal{D} , which are essentially bounded from above and from below by the constants ρ_{\min} and ρ_{\max} , respectively.

In this section we comment on why an analogous result could not be shown so far either for PBSOPs involving the TE band structure of a two-dimensional photonic crystal, or for PBSOPs involving the band structure of a three-dimensional photonic crystal.

In the following, we assume that Ω is a primitive domain of a Bravais lattice $\Lambda \subset \mathbb{R}^2$ of rank 2, and the \mathbb{B} is the first Brillouin zone of Λ . Recall that $\lambda_j^{\text{TE}}(\rho, \mathbf{k})$ denotes the j -th smallest eigenvalue of the eigenvalue equation (4.94) on page 87 for every coefficient $\rho \in \mathcal{D}$ and every vector $\mathbf{k} \in \mathbb{B}$. If the coefficient ρ belongs to the admissible set \mathcal{C} , which was defined in (5.3), then it represents the Λ -periodic relative electric permittivity function ε_r of a two-dimensional photonic crystal according to $\rho = 1/\varepsilon_r|_{\Omega}$.

Given some fixed number $n \in \mathbb{N}$, we shall call a functional $J^{\text{TE}} : \mathcal{D} \rightarrow \mathbb{R}$ a *TE goal functional* if it is of the form

$$J^{\text{TE}}(\rho) = \Upsilon^{\text{TE}}(\omega_1^{\text{TE}}(\rho, \cdot), \dots, \omega_n^{\text{TE}}(\rho, \cdot)) \quad \text{for all } \rho \in \mathcal{D}, \quad (5.8)$$

where

$$\Upsilon^{\text{TE}} : \underbrace{C^0(\mathbb{B}, \mathbb{R}) \times \dots \times C^0(\mathbb{B}, \mathbb{R})}_{n\text{-times}} \rightarrow \mathbb{R},$$

is a locally Lipschitz continuous functional. Supposing that a minimizer $\rho_* \in \mathcal{C}$ of the TE goal functional J^{TE} represents the relative electric permittivity function of a two-dimensional photonic crystal, whose TE band structure exhibits some desired property, we consider the following scheme of optimization problems.

Problem 5.13. *Given a TE goal functional $J^{\text{TE}} : \mathcal{D} \rightarrow \mathbb{R}$ of the form (5.8), find an density function $\rho_* \in \mathcal{C}$, such that*

$$J^{\text{TE}}(\rho_*) \leq J^{\text{TE}}(\rho) \quad \text{for all } \rho \in \mathcal{C}.$$

One can show that every TE goal functional is locally Lipschitz continuous on the extended admissible set $\bar{\mathcal{C}}$, which was defined by (5.5). The question is, whether or not one can show the existence of optimal solutions in $\bar{\mathcal{C}}$, as for Problem 5.13.

Recall that the existence of optimal solutions of Problem 5.13 followed from the weak-*compactness of the extended admissible set $\bar{\mathcal{C}}$ and the weak-*continuity of the TM eigenvalue functions $\lambda_j^{\text{TM}}(\cdot, \mathbf{k})$. The latter result was established by Theorem 5.11 in the previous section. The proof of Theorem 5.11 relies on the result stated by Lemma 5.10, namely that for arbitrary sequences $\{\rho^{(l)}\}_{l \in \mathbb{N}}$, $\{f^{(l)}\}_{l \in \mathbb{N}}$ and $\{g^{(l)}\}_{l \in \mathbb{N}}$ in $\bar{\mathcal{C}}$, Z , and Z , the convergence of $m^{\text{TM}}(\rho^{(l)})(f^{(l)}, g^{(l)})$ to $m^{\text{TM}}(\rho)(f, g)$ as l tends to infinity holds, if the sequence $\{\rho^{(l)}\}_{l \in \mathbb{N}}$ weak-*converges in \mathcal{E} to $\rho \in \bar{\mathcal{C}}$, and if the sequences $\{f^{(l)}\}_{l \in \mathbb{N}}$ and $\{g^{(l)}\}_{l \in \mathbb{N}}$ converge strongly in Z to $f \in Z$ and $g \in Z$, respectively.

When trying to adapt the proof of Theorem 5.11 to setting of Problem 5.13, one finds that the following statement needs to be true: Given a sequence $\{\rho^{(l)}\}_{l \in \mathbb{N}}$ in $\bar{\mathcal{C}}$ that weak-*converges in \mathcal{E} to a density function $\rho \in \bar{\mathcal{C}}$, as well as two sequences $\{u^{(l)}\}_{l \in \mathbb{N}}$ and $\{v^{(l)}\}_{l \in \mathbb{N}}$ in V that converge strongly in Z and weakly in V to $u \in V$ and $v \in V$, respectively, it follows that $a_{\mathbf{k}}^{\text{TE}}(\rho^{(l)})(u^{(l)}, v^{(l)})$ converges to $a_{\mathbf{k}}^{\text{TE}}(\rho)(u, v)$ as l tends to infinity. Unfortunately, this statement cannot be shown. In fact, for $\mathbf{k} = \mathbf{0}$ the statement is known to be false. In the following we shall make this notion more precise.

Suppose that $\{u^{(l)}\}_{l \in \mathbb{N}}$ is a sequence in V that converges weakly in V to a function $u \in V$. Then, the sequence of weakly defined gradients $\{\nabla u^{(l)}\}_{l \in \mathbb{N}}$ converges weakly in Z (see e.g. Example 6.4(3) in [4]). Let $\{\rho^{(l)}\}_{l \in \mathbb{N}}$ be a sequence in $\bar{\mathcal{C}}$ that weak-*converges in \mathcal{E} to a function $\rho \in \bar{\mathcal{C}}$. Then, the sequence $\{\rho^{(l)} \nabla u^{(l)}\}_{l \in \mathbb{N}}$ does not converge weakly in Z to $\rho \nabla u$, in general. A counter-example is constructed, for example, in Section 5.1 in [21]. Hence, we cannot expect $a_{\mathbf{0}}^{\text{TE}}(\rho^{(l)})(u^{(l)}, v)$ to converge to $a_{\mathbf{0}}^{\text{TE}}(\rho)(u, v)$ for every $v \in V$.

Sequences of the form $\{\rho^{(l)} \nabla u^{(l)}\}_{l \in \mathbb{N}}$ typically occur in problems, which are studied in homogenization theory. An important concept that is related to these sequences is the concept of *H-convergence* (see [56]), which we shall briefly introduce. In the following, we denote by $D \subset \mathbb{R}^r$ an open, simply connected, and bounded domain, where $r \in \mathbb{N}$ is an arbitrary space dimension. Then, we defined and for every two constants $\alpha_{\min}, \alpha_{\max} \in \mathbb{R}$, that satisfy $0 < \alpha_{\min} < \alpha_{\max}$, the set

$$\begin{aligned} M(D, \alpha_{\min}, \alpha_{\max}) := \{ \mathbf{A} \in L^\infty(D, \mathbb{R})^{r \times r} \mid & \langle \mathbf{A}(\cdot) \mathbf{x}, \mathbf{x} \rangle \geq \alpha_{\min} |\mathbf{x}|^2, \\ & |\mathbf{A}(\cdot) \mathbf{x}| \leq \alpha_{\max} |\mathbf{x}| \\ & \text{for all } \mathbf{x} \in \mathbb{R}^r \text{ a.e. on } D \}. \end{aligned} \quad (5.9)$$

Clearly, the set $M(D, \alpha_{\min}, \alpha_{\max})$ consists of specific matrix-valued functions defined on D . A sequence $\{\mathbf{A}^{(l)}\}_{l \in \mathbb{N}}$ in $M(D, \alpha_{\min}, \alpha_{\max})$ is said to *H-converge* to a matrix-valued function $\mathbf{A} \in M(D, \beta_{\min}, \beta_{\max})$, if for every continuous, linear functional $f \in H^{-1}(D)$ the solution $u^{(l)} \in H_0^1(D)$ of the weakly formulated boundary value problem

$$\int_D \mathbf{A}^{(l)} \overline{\nabla u^{(l)}} \cdot \nabla v = f(v) \quad \text{for all } v \in H_0^1(D),$$

converges weakly in $H_0^1(D)$ to the solution $u \in H_0^1(D)$ of the weakly formulated boundary value problem

$$\int_D \mathbf{A} \overline{\nabla u} \cdot \nabla v = f(v) \quad \text{for all } v \in H_0^1(D),$$

and if $\mathbf{A}^{(l)} \nabla u^{(l)}$ converges weakly in $L^2(D)$ to the vector field $\mathbf{A} \nabla u$ as l tends to ∞ . It is known that for every sequence $\{\mathbf{S}^{(l)}\}_{l \in \mathbb{N}}$ of symmetric-matrix-valued

functions in $M(D, \alpha_{\min}, \alpha_{\max})$ there exists a symmetric-matrix-valued function $\mathbf{S} \in M(D, \alpha_{\min}, \alpha_{\max})$, such that a subsequence of $\{\mathbf{S}^{(l)}\}_{l \in \mathbb{N}}$ H-converges to \mathbf{S} (see Theorem 13.2(iii) and Proposition 13.6 in [21]). This notion is expressed in stating that the set $M_2(\alpha_{\min}, \alpha_{\max})$ is *G-compact* (see Theorem 13.2 in [21]).

By definition, the concepts of H-convergence and G-compactness are linked to weakly formulated boundary value problems of the form

$$\int_D \mathbf{A} \overline{\nabla u} \cdot \nabla v = f(v) \quad \text{for all } v \in H_0^1(\Omega),$$

where $u \in H_0^1(D)$ is the unknown function, where $f \in H^{-1}(D)$ is a given right-hand side, and where $\mathbf{A} \in M(D, \rho_{\min}, \rho_{\max})$ is a matrix-valued coefficient. Unfortunately, it seems that no attempts have been made so far to generalize these concepts. Anyway, we were not able to find conclusive evidence in the literature as to whether or not the results related to H-convergence also apply to weakly formulated boundary value problems of the form

$$\int_{\Omega} \mathbf{A} \overline{(\nabla + \mathbf{i}\mathbf{k})u} \cdot (\nabla + \mathbf{i}\mathbf{k})v = f(v) \quad \text{for all } v \in V,$$

where $u \in V$ is the unknown function, $f \in V^*$ is a prescribed right-hand side, $\mathbf{k} \in \mathbb{B}$ is a given vector, and $\mathbf{A} \in M(\Omega, \rho_{\min}, \rho_{\max})$ is a given matrix-valued coefficient. Provided that the set $M(\Omega, \rho_{\min}, \rho_{\max})$ is also G-compact with respect to these problems, one could show that every TE goal functional attains a minimum on the set $M(\Omega, \rho_{\min}, \rho_{\max})$. The corresponding minimizer would then be a symmetric-matrix-valued function, representing the relative electric permittivity function of an anisotropic medium.

The situation is very similar for PBSOPs that involve band structures of three-dimensional photonic crystals. Henceforth, we assume that $\Lambda \subset \mathbb{R}^3$ is a Bravais lattice of rank 3 with primitive cell Ω and first Brillouin zone \mathbb{B} . We assume that the admissible set \mathcal{C} and the extended admissible set $\overline{\mathcal{C}}$ are defined according to (5.2) and (5.4).

Under the assumptions of Section 5.1, one is not able to prove the existence of optimal solutions of Problem 5.1 in the extended admissible set $\overline{\mathcal{C}}$ by adapting the proof of Theorem 5.11. What is missing, again, is a suitable analogue of Lemma 5.10. More precisely, given a sequence $\{\rho^{(l)}\}_{l \in \mathbb{N}}$ in $\overline{\mathcal{C}}$ that weak-*converges in \mathcal{E} to $\rho \in \overline{\mathcal{C}}$, and two sequences $\{\mathbf{u}^{(l)}\}_{l \in \mathbb{N}}$ and $\{\mathbf{v}^{(l)}\}_{l \in \mathbb{N}}$ in $\mathbf{V}_{\mathbf{k}}$ that converge weakly in $\mathbf{V}_{\mathbf{k}}$ and strongly in \mathbf{Z} to $\mathbf{u} \in \mathbf{V}_{\mathbf{k}}$ and $\mathbf{v} \in \mathbf{V}_{\mathbf{k}}$, respectively, one cannot conclude that $a_{\mathbf{k}}(\rho^{(l)})(\mathbf{u}^{(l)}, \mathbf{v}^{(l)})$ converges to $a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{v})$ as l tends to infinity. As in the TE setting, one could prove the existence of optimal solutions in the set $M(\Omega, \rho_{\min}, \rho_{\max})$, if the results related to H-convergence also apply to weakly

formulated boundary value problems of the form

$$\int_{\Omega} \mathbf{A}(\nabla + i\mathbf{k}) \times \mathbf{u} \cdot (\nabla + i\mathbf{k})\mathbf{v} = f(\mathbf{v}) \quad \text{for all } v \in \mathbf{V}_{\mathbf{k}},$$

where $\mathbf{u} \in \mathbf{V}_{\mathbf{k}}$ is the unknown function, $f \in \mathbf{V}_{\mathbf{k}}^*$ is a prescribed right-hand side, $\mathbf{k} \in \mathbb{B}$ a given vector and $\mathbf{A} \in M(\Omega, \rho_{\min}, \rho_{\max})$ a given coefficient. An optimal solution in $M(\Omega, \rho_{\min}, \rho_{\max})$ would then be a symmetric-matrix-valued function representing the relative electric permittivity function of an anisotropic medium. To the best of our knowledge, however, the question as to whether or not the concepts of H-convergence and G-compactness also apply to the above type of boundary value problem has not been considered in the literature so far.

5.5 Optimization Goals

In the previous sections we studied general aspects of photonic band structure optimization problems (PBSOPs). In this section we discuss several concrete optimization problems and introduce a specific class of model problems, which will be further investigated in this dissertation.

The most common PBSOPs, which are studied in the literature, are so-called *photonic band gap maximization problems* (PBGMPs). A photonic band gap maximization problem consists in finding a medium structure of a photonic crystal, which yields a band structure that exhibits a band gap as large as possible. The main reason for the popularity of PBGMPs is, that photonic crystals, whose band structures exhibit large band gaps, are expected to play an important role in the creation of certain nano-optical devices. Recall that photonic band gaps correspond to ranges of optical frequencies, at which the photonic crystal inhibits any light propagation. Therefore, photonic crystals with large band gaps are favorable for the design of mirrors or nano-scale optical wave guides.

In the following we present two PBSOPs, which were discussed in the literature and which fit into the problem scheme given by Problem 5.1 on page 92. It is easy to see that every PBSOP, which fits into that scheme is uniquely identified by the respective goal functional J .

In their work on the maximization of photonic band gaps in two-dimensional photonic crystals [29], the authors Cox and Dobson considered TM goal functionals $J_{j;\omega_0}^{\text{TM}} : \mathcal{D} \rightarrow \mathbb{R}$, which were defined by

$$J_{j;\omega_0}^{\text{TM}}(\rho) := - \min_{\mathbf{k} \in \mathbb{B}} \min \{ \omega_{j+1}^{\text{TM}}(\rho, \mathbf{k})^2 - \omega_0^2, \omega_0^2 - \omega_j^{\text{TM}}(\rho, \mathbf{k})^2 \} \quad (5.10)$$

for all $\rho \in \mathcal{D}$. Here, $j \in \mathbb{N}$ is a given index, and ω_0 is a given real number. Under the assumption that there exists an initial density function $\rho^{(0)} \in \overline{\mathcal{C}}$, such that the

photonic crystal represented by $\rho^{(0)}$ exhibits a band gap between the j -th and the $(j + 1)$ -th TM band, the number ω_0 can be chosen such that

$$\max_{\mathbf{k} \in \mathbb{B}} \omega_j^{\text{TM}}(\rho^{(0)}, \mathbf{k}) < \omega_0 < \min_{\mathbf{k} \in \mathbb{B}} \omega_{j+1}^{\text{TM}}(\rho^{(0)}, \mathbf{k}).$$

Such a number ω_0 represents a so-called *gap frequency* in the TM band structure of the photonic crystal, which is represented by $\rho^{(0)}$. The minimization problem for the corresponding goal functional $J_{j;\omega}^{\text{TM}}$ aims at increasing the minimal distance between the j -th and the $(j + 1)$ -th TM band to the gap frequency ω_0 . By this, Cox and Dobson were able to widen the TM band gaps of several two-dimensional photonic crystals.

In a second paper on band structure optimization of two-dimensional photonic crystals [30], Cox and Dobson used another type of TE goal functionals. These functionals $J_{j;\omega}^{\text{TE}} : \mathcal{D} \rightarrow \mathbb{R}$ are given by

$$J_{j;\omega}^{\text{TE}}(\rho) := - \min_{\mathbf{k} \in \mathbb{B}} \min \{ \omega_{j+1}^{\text{TE}}(\rho, \mathbf{k})^2 - \omega(\mathbf{k})^2, \omega(\mathbf{k})^2 - \omega_j^{\text{TE}}(\rho, \mathbf{k})^2 \} \quad (5.11)$$

for all $\rho \in \mathcal{D}$. Here, $\omega : \mathbb{B} \rightarrow \mathbb{R}$ denotes a prescribed, continuous function. Under the assumption that there exists an index $j \in \mathbb{N}$ and an initial density function $\rho^{(0)} \in \mathcal{C}^{\text{TE}}$, such that the corresponding j -th and $(j + 1)$ -th TE band are strictly separated (see (4.75) on page 80), the function ω can be chosen such that

$$\omega_j^{\text{TE}}(\rho^{(0)}, \mathbf{k}) < \omega(\mathbf{k}) < \omega_{j+1}^{\text{TE}}(\rho^{(0)}, \mathbf{k}) \quad \text{for all } \mathbf{k} \in \mathbb{B}. \quad (5.12)$$

The minimization problem for the corresponding goal functional $J_{j;\omega}^{\text{TE}}$ aims at increasing the minimal distance of the j -th and the $(j + 1)$ -th TE band to the graph of the function ω . In an iterative optimization algorithm, one can modify the function ω slightly after each iteration and thus create a sequence of functions $\{\omega_l\}_{l \in \mathbb{N}}$, such that every element of that sequence satisfies (5.12). Cox and Dobson devised an iterative algorithm, in which the sequence $\{\omega_l\}_{l \in \mathbb{N}}$ converged to a constant function. By this, they were able to open band gaps in the TE band structures of some two-dimensional photonic crystals.

In the following sections, we consider PBSOPs for three-dimensional photonic crystals, which are given by goal functionals $J_j : \mathcal{D} \rightarrow \mathbb{R}$ of the form

$$J_j(\rho) := \max_{\mathbf{k} \in \mathbb{B}} \lambda_j(\rho, \mathbf{k}) - \min_{\mathbf{k} \in \mathbb{B}} \lambda_{j+1}(\rho, \mathbf{k}) \quad \text{for all } \rho \in \mathcal{D}, \quad (5.13)$$

for some index $j \in \mathbb{N}$. We will refer to these functions as *gap width functionals*, since minimization problems for these functionals aim at widening the gap width of a band gap between the j -th and the $(j + 1)$ -th photonic band. Given an index $j \in \mathbb{N}$, an optimization algorithm for the goal functional J_j can also be used to open up a band gap between the j -th and the $(j + 1)$ -th photonic band, provided

that these bands are strictly separated. In Chapter 9 we present some numerical examples, where band gaps between separated bands could be opened.

For completeness we also define for every index $j \in \mathbb{N}$ the TM and TE analogues $J_j^{\text{TM}} : \mathcal{D} \rightarrow \mathbb{R}$ and $J_j^{\text{TE}} : \mathcal{D} \rightarrow \mathbb{R}$ of J_j by

$$J_j^{\text{TM}}(\rho) := \max_{\mathbf{k} \in \mathbb{B}} \lambda_j^{\text{TM}}(\rho, \mathbf{k}) - \min_{\mathbf{k} \in \mathbb{B}} \lambda_{j+1}^{\text{TM}}(\rho, \mathbf{k}) \quad \text{for all } \rho \in \mathcal{D}, \quad (5.14)$$

$$J_j^{\text{TE}}(\rho) := \max_{\mathbf{k} \in \mathbb{B}} \lambda_j^{\text{TE}}(\rho, \mathbf{k}) - \min_{\mathbf{k} \in \mathbb{B}} \lambda_{j+1}^{\text{TE}}(\rho, \mathbf{k}) \quad \text{for all } \rho \in \mathcal{D}. \quad (5.15)$$

These goal functionals were also used by Kao, Osher and Yablonovitch in their work on the maximization of band gaps in two-dimensional photonic crystals (cf. [44]).

We compared the goal functionals $J_{j;\omega_0}^{\text{TM}}$ and J_j^{TM} by conducting some numerical experiments. In these experiments we found that solving a PBSOPs for the TM gap width functionals J_j^{TM} usually yields larger band gaps than solving such problems for the goal functionals $J_{j;\omega_0}^{\text{TM}}$. This is mostly due to the fact, that an optimization algorithm for the TM gap width functionals J_j^{TM} can freely change the absolute position of the band gap in order to widen it. In contrast to this, optimization algorithms for the goal functionals $J_{j;\omega_0}^{\text{TM}}$ are restricted in their ability to change the position of the band gap by the fixed gap frequency ω_0 .

For some applications, however, prescribing a fixed gap frequency ω_0 is essential. Consider, for example, the problem of designing a photonic crystal that should inhibit light propagation around a prescribed center frequency. The center frequency could be determined, for example, by a given light source, such as a specific laser device. For this problem it is more adequate, of course, to choose a goal functional similar to $J_{j;\omega_0}^{\text{TM}}$ for some index $j \in \mathbb{N}$, and to use to the prescribed center frequency as ω_0 . Another situation, where it can be useful to prescribe a certain gap frequency ω_0 , is when one is interested in finding two-dimensional photonic crystals, which exhibit maximal complete band gaps. As was stated in Section 4.7 a complete band gap consists of all frequencies lying in a TM band gap as well as in a TE band gap. Maximizing complete band gaps thus amounts to maximizing TM and TE band gaps, which contain a common frequency ω_0 .

In this work, we are mostly interested in finding photonic crystals with band gaps that are as large as possible. Therefore, we shall only consider PBSOPs for the gap width functionals J_j , J_j^{TM} , and J_j^{TE} in the following chapters.

Chapter 6

Nonsmooth Analysis

In Section 5.2 we showed that the goal functionals, which arise in photonic band structure optimization problems (PBSOPs), are typically locally Lipschitz continuous. It was also emphasized that the goal functionals in general fail to be differentiable. However, in order to devise a local optimization algorithm for PBSOPs one needs some information about the local monotonicity of the goal functionals. For locally Lipschitz continuous functionals, this information can be obtained from so-called generalized differentials. Generalized differentials are a fundamental tool of nonsmooth analysis and were first introduced by Clarke (see e.g. [22]). In this chapter we list some important results concerning generalized differentials. In Section 6.1 we define generalized directional derivatives, as well as generalized differentials and discuss some of their properties. In Section 6.2 we provided some rules of calculus for generalized differentials. These rules are used in Section 6.3 in order to determine the generalized differentials of the gap width functionals for three-dimensional photonic crystals, which were introduced in Section 5.5. In Section 6.4 we give a characterization of the gap width functionals for two-dimensional photonic crystals.

6.1 Generalized Differentials

Throughout this section we denote by X a real Banach space, which is endowed with a norm $\|\cdot\|_X$. By X^* we denote its normed dual. Given a functional $f : X \rightarrow \mathbb{R}$, a point $x \in X$, and a vector $v \in X$, we define after Clarke (see Chapter 2.1 in [22]) the upper limit

$$f^\circ(x; v) := \limsup_{\substack{y \rightarrow x \\ t \rightarrow 0^+}} \frac{f(y + tv) - f(y)}{t}.$$

For locally Lipschitz continuous functionals f this upper limit is a real number, which is called the *generalized directional derivative of f at x in direction v* . For functionals, that are continuously differentiable in a neighbourhood of x , $f^\circ(x; v)$ coincides with the one-sided directional derivative of f at x in direction v . The following proposition summarizes important properties of generalized directional derivatives. The proposition corresponds to Proposition 1.1 in Chapter 2 in [23].

Proposition 6.1. *Let $f : X \rightarrow \mathbb{R}$ be a functional, which is c -Lipschitz continuous near a point $x \in X$ for some local Lipschitz constant $c > 0$. Then, the following statements are valid.*

(a) $f^\circ(x; \cdot) : X \rightarrow \mathbb{R}$ is positively homogeneous, subadditive, and satisfies

$$|f^\circ(x; v)| \leq c\|v\|_X \quad \text{for all } v \in X.$$

(b) $f^\circ(x; \cdot) : X \rightarrow \mathbb{R}$ is c -Lipschitz continuous near x .

(c) $f^\circ : X \times X \rightarrow \mathbb{R}$ is upper semicontinuous.

(d) $f^\circ(x; -v) = (-f)^\circ(x; v)$ for all $v \in X$.

To a certain extent the sign of the functional $f^\circ(x; \cdot) : X \rightarrow \mathbb{R}$ is determined by the local monotonicity of the functional f at the point x . This notion is made precise by the following proposition. We denote by $B_r(x)$ the open ball in X with radius $r > 0$ and center $x \in X$.

Proposition 6.2. *Let $f : X \rightarrow \mathbb{R}$ be a locally Lipschitz continuous functional, let $x \in X$ be a point, and let $V \subset X$ be a closed, convex cone with vertex at the origin.*

(a) *If f is monotonically increasing at x in all directions of V , i.e., if there exists a radius $r > 0$, such that*

$$f(x + w) \geq f(x) \quad \text{for all } w \in V \cap B_r(0),$$

then $f^\circ(x; v) \geq 0$ for all $v \in V$.

(b) *If f is monotonically increasing in an open neighbourhood U of x in all directions of V , i.e., if there exists a radius $r > 0$, such that*

$$f(y + w) \geq f(y) \quad \text{for all } y \in U, w \in V \cap B_r(0),$$

then $f^\circ(x; v) \geq 0$ and $f^\circ(x; -v) \leq 0$ for all $v \in V$.

Proof. Let $v \in C$. Then, there exists a number $t_0 > 0$, such that $tv \in B_r(0)$ for all $t \in [0, t_0]$. Hence,

$$\begin{aligned} f^\circ(x; v) &= \limsup_{\substack{y \rightarrow x \\ t \rightarrow 0}} \frac{f(y + tv) - f(y)}{t} \\ &\geq \limsup_{t \rightarrow 0^+} \frac{f(x + tv) - f(y)}{t} \\ &= \limsup_{t \rightarrow 0, t \in [0, t_0]} \frac{f(x + tv) - f(y)}{t} \\ &\geq 0, \end{aligned}$$

which proves the inequality $f^\circ(x; v) \geq 0$ in (a) and (b). Furthermore, we have

$$\begin{aligned} f^\circ(x; -v) &= (-f)^\circ(x; v) \\ &= \limsup_{\substack{y \rightarrow x \\ t \rightarrow 0}} \frac{f(y) - f(y + tv)}{t} \\ &= \limsup_{\substack{y \rightarrow x, y \in U \\ t \rightarrow 0, t \in [0, t_0]}} \frac{f(y) - f(y + tv)}{t} \\ &\leq 0, \end{aligned}$$

which establishes the inequality $f^\circ(x; -v) \leq 0$ in (b). \square

In contrast to one-sided directional derivatives, the sign of the generalized directional derivative does not provide conclusive information about the local monotonicity of a locally Lipschitz continuous functional. We illustrate this by the following example.

Example 6.3. Consider the function $f : \mathbb{R} \rightarrow \mathbb{R}$, defined as $f(x) := \sqrt{1 - |x|} - 1$ for all $x \in \mathbb{R}$ (see Figure 6.1). This function is Lipschitz continuous on \mathbb{R} , the minimal Lipschitz constant being $1/2$. Moreover, f is continuously differentiable at all points except at the origin. Explicit computation shows that the generalized directional derivative at the origin is given by $f^\circ(0; v) = |v|/2$. Hence, $f^\circ(0; v) > 0$ for all $v \in \mathbb{R}$, which intuitively corresponds well to the fact that f increases monotonically at the origin in all directions. However, by Proposition 6.1(c) we also have $(-f)^\circ(0; v) = f^\circ(0; v) > 0$ for all $v \in \mathbb{R}$, even though the function $-f$ decreases monotonically at the origin in all directions.

By Proposition 6.1(a) the generalized directional derivative of a locally Lipschitz continuous functional constitutes a positively homogeneous, subadditive, bounded

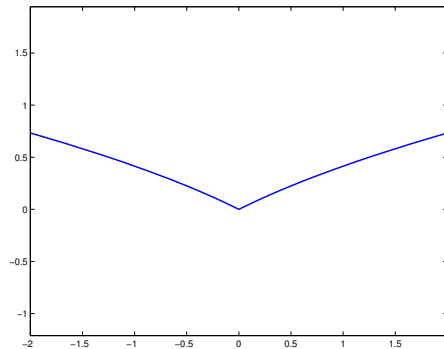


Figure 6.1: Graph of the function f in Example 6.3 at the origin.

functional on X . According to the Hahn–Banach Theorem, the generalized directional derivative thus majorizes at least one bounded, linear functional on X . Given a locally Lipschitz functional $f : X \rightarrow \mathbb{R}$ and a point $x \in X$, we therefore have that the set

$$\partial f(x) := \{g \in X^* \mid g(v) \leq f^\circ(x; v) \text{ for all } v \in X\}. \quad (6.1)$$

is non-empty. This set $\partial f(x)$ is called the *generalized differential of f at x* . The following proposition summarizes important properties of generalized differentials. The assertions of the proposition correspond to those stated in the Propositions 2.1.2 and 2.3.1 in [22].

Proposition 6.4. *Let $f : X \rightarrow \mathbb{R}$ be L -Lipschitz continuous near $x \in X$. Then,*

- (a) $\partial f(x)$ is a non-empty, convex, weak*-compact, and bounded subset of X^* .
In particular,

$$\|p\|_{X^*} \leq L \quad \text{for all } p \in \partial f(x).$$

- (b) $f^\circ(x; v) = \max\{p(v) \mid p \in \partial f(x)\}$ for all $v \in X$.

- (c) $\partial(-f)(x) = -\partial f(x)$.

As the name already suggests, generalized differentials are generalizations of classical differentials or derivatives for locally Lipschitz continuous functionals. In particular, they provide necessary conditions for the existence of local extrema according to the following lemma.

Lemma 6.5. *If a locally Lipschitz continuous functional $f : X \rightarrow \mathbb{R}$ attains a local extremum at some point $x \in X$, then $0 \in \partial f(x)$.*

For a proof of Lemma 6.5, see Proposition 2.3.2 in [22]. There also exists an analog of the Mean-Value Theorem for generalized differentials, which is stated in the following Theorem of Lebourg (cf. [52]).

Theorem 6.6 (Lebourg). *Let $f : X \rightarrow \mathbb{R}$ be Lipschitz continuous on a convex, open subset K of X . Then, for every two points $x, y \in K$ there exists a point $z \in \text{conv}\{x, y\}$ and a continuous, linear functional $p \in \partial f(z)$, such that*

$$f(y) - f(x) = p(y - x).$$

The following proposition reveals the relationship between generalized differentials and other concepts of differentiability. The assertions of the proof are established in Section 2.2 in [22].

Proposition 6.7. *Let $f : X \rightarrow \mathbb{R}$ be a locally Lipschitz continuous functional, and let U be an open neighbourhood of some point $x \in X$.*

- (a) *If f is convex on U , then $\partial f(x)$ coincides with the subdifferential of f at x .*
- (b) *If f is Gâteaux differentiable or Fréchet differentiable at x with derivative $Df(x)$, then $Df(x) \in \partial f(x)$.*
- (c) *If f is continuously Gâteaux differentiable on U with derivative $Df(x)$ at x , then $\{Df(x)\} = \partial f(x)$.*
- (d) *The functional f is strictly differentiable in U with strict derivative $D_s f(x)$ at x , if and only if $\{D_s f(x)\} = \partial f(x)$.*

Recall that the *subdifferential* of a convex functional $f : X \rightarrow \mathbb{R}$ at a point $x \in X$ is defined as set of all continuous linear functionals $p \in X^*$, for which

$$f(y) - f(x) \geq p(y - x) \quad \text{for all } y \in X.$$

In the literature the subdifferential of f at a point x is also commonly denoted by $\partial f(x)$.

The concept of strict differentiability is less well-known, which is why we briefly comment on it here. A functional $f : X \rightarrow \mathbb{R}$ is called *strictly differentiable* at a point $x \in X$, if there exists a continuous linear functional $p \in X^*$, such that

$$\lim_{\substack{y \rightarrow x \\ t \rightarrow 0+}} \frac{f(y + tv) - f(y)}{t} = p(v) \quad \text{for all } v \in X.$$

The functional p is then called the *strict differential of f at x* and denoted by $D_s f(x)$. It can be shown that continuous Gâteaux differentiability at a point $x \in X$

implies strict differentiability at that point (see Corollary to Proposition 2.2.1 in [22]). Moreover, it is easy to see that strict differentiability implies Gâteaux differentiability.

It should be noted that Clarke, as well as other authors, call the set $\partial f(x)$ the “generalized gradient” of f at x . We prefer the term “generalized differential”, however, since it alludes to the fact, that $\partial f(x)$ is a subset of X^* . This corresponds to the notion of the differentials of a Gâteaux differentiable functional $f : X \rightarrow \mathbb{R}$. By definition, the differential $Df(x)$ of such a functional at a given point $x \in X$ is an element of X^* . Now suppose that Y is another real Banach space, and that $\langle \cdot, \cdot \rangle : Y \times X \rightarrow \mathbb{R}$ is a bilinear form, such that $(Y, X, \langle \cdot, \cdot \rangle)$ is a *dual pair* in the sense of functional analysis (see e.g. Definition VIII.3.1 in [75]). If there exists a uniquely defined element $u \in Y$, such that $D\varphi(x) = \langle u, \cdot \rangle$ then u is called the *gradient of f at x* with respect to $\langle \cdot, \cdot \rangle$. If X is a Hilbert space, one usually chooses the inner product on X as the bilinear form $\langle \cdot, \cdot \rangle$. The gradient of f at a point x is then commonly denoted by $\nabla f(x)$. The connection between differentials and gradients in a Hilbert space is established by the Rieszian isomorphism $R : X \rightarrow X^*$ through $\nabla f(x) = R^{-1}(Df(x))$.

According to the discussion above, it is quite natural to define generalized gradients as follows. Let X and Y be real Banach spaces, and let $\langle \cdot, \cdot \rangle : Y \times X \rightarrow \mathbb{R}$ be a bilinear form, such that $(Y, X, \langle \cdot, \cdot \rangle)$ is a dual pair. Given a locally Lipschitz continuous functional $f : X \rightarrow \mathbb{R}$ and a point $x \in X$, the set

$$\{u \in Y \mid \langle u, \cdot \rangle \in \partial f(x)\} \quad (6.2)$$

is called the *generalized gradient of f at x* with respect to $\langle \cdot, \cdot \rangle$. If X is a Hilbert space coinciding with Y , and if $\langle \cdot, \cdot \rangle$ is the inner product on X , then the set in (6.2) is often identified with the generalized differential of f at x and also denoted by $\partial f(x)$.

We conclude this section with the following example.

Example 6.8. Consider again the function f defined in Example 6.3. The function is defined on \mathbb{R} . The dual space \mathbb{R}^* of \mathbb{R} consists of all functionals $p : \mathbb{R} \rightarrow \mathbb{R}$, which are of the form $p(v) = \alpha v$, where $\alpha \in \mathbb{R}$. By (6.1) a functional $p \in \mathbb{R}^*$ belongs to $\partial f(0)$, if and only if $p(v) \leq f^\circ(x, v)$ for all $v \in \mathbb{R}$. Hence, one easily deduces that

$$\partial f(0) = \left\{ p \in \mathbb{R}^* \mid p(v) = \alpha v, \alpha \in \left[-\frac{1}{2}, \frac{1}{2} \right] \right\}.$$

In view of Lemma 6.5 it is not surprising, that the zero function on \mathbb{R} is contained in $\partial f(0)$, since f attains a local minimum at 0. Since \mathbb{R} is a Hilbert space, one can also define a generalized gradient of f at x as a subset of \mathbb{R} . Clearly, this subset is given by the interval $[-1/2, 1/2]$.

6.2 Generalized Differential Calculus

In this section we list some rules of calculus for generalized differentiation. As we shall see below, most of these rules only establish certain inclusion relations. Therefore, these rules can only be used to determine supersets of generalized differentials. Still, these rules will prove useful in the following section, where we aim to characterize the generalized differentials of the gap width functionals introduced in Section 5.5.

As in the previous section, we denote by X a real Banach space. Our first rule of calculus, which corresponds to Proposition 2.3.1 in [22], states that generalized differentiation is a homogeneous operation.

Proposition 6.9. *Let $f : X \rightarrow \mathbb{R}$ be a functional, which is Lipschitz continuous near a point $x \in X$, and let $\theta \in \mathbb{R}$ be a real number. Then we have*

$$\partial(\theta f)(x) = \theta \partial f(x).$$

The next proposition, which corresponds to Proposition 2.3.3 in [22], states that generalized differentials are subadditive in the sense of set inclusions.

Proposition 6.10. *Let $f_1, \dots, f_N : X \rightarrow \mathbb{R}$ be a finite number of functionals for some $N \in \mathbb{N}$, which are Lipschitz continuous near a point $x \in X$. Then, we have*

$$\partial(f_1 + \dots + f_n)(x) \subseteq \partial f_1(x) + \dots + \partial f_N(x).$$

The following theorem establishes a chain rule for certain continuously differentiable functionals with locally Lipschitz continuous functionals. This theorem is a particularization of Theorem 2.3.9 in [22].

Theorem 6.11. *Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a continuously differentiable function, and let $g : X \rightarrow \mathbb{R}$ be a functional, which is Lipschitz continuous near a point $x \in X$. Then, the functional $f \circ g$ is Lipschitz continuous near x , and we have*

$$\partial(f \circ g)(x) = f'(g(x))\partial g(x) := \{f'(g(x))q \mid q \in \partial g(x)\},$$

where $f'(g(x)) \in \mathbb{R}$ denotes the derivative of f at the point $g(x)$.

The following theorem is a particularization of Theorem 2.3.10 in [22]. It establishes another chain rule for generalized differentials of functionals, which are given by compositions of locally Lipschitz continuous functions and continuously Gâteaux differentiable operators. Just like Proposition 6.10, the theorem only establishes an inclusion relation.

Theorem 6.12. *Let Y be a second real Banach space, and let $G : X \rightarrow Y$ be an operator, which is continuously Gâteaux differentiable in a neighbourhood of a point $x \in X$. Furthermore, let $f : Y \rightarrow \mathbb{R}$ be a functional, which is Lipschitz continuous near $G(x)$. Then, the functional $f \circ G$ is Lipschitz continuous near x , and we have*

$$\partial(f \circ G)(x) \subseteq \{p \circ DG(x) \mid p \in \partial f(G(x))\},$$

where $DG(x)$ denotes the Gâteaux derivative of G at x .

Before we can state the next rule of calculus for generalized differentiation, we need to introduce some concepts from measure theory. The concepts we present here are discussed in more detail in [67]. Suppose that T is a topological Hausdorff space and that \mathcal{T} is a σ -algebra on T containing the open subsets of T . Then, a measure μ on \mathcal{T} is called *locally finite*, if for every point $t \in T$ there exists an open neighbourhood U of t , such that $\mu(U) < \infty$. The measure μ is called *inner regular*, if

$$\mu(S) = \sup\{\mu(K) \mid K \subset S, K \text{ compact}\} \quad \text{for all } S \in \mathcal{T}.$$

Let $\mathcal{B}(T)$ be the Borel σ -algebra on T . Recall that $\mathcal{B}(T)$ is defined as the smallest σ -algebra on T containing all open subsets of T . It follows from the definition of σ -algebras that $\mathcal{B}(T)$ is equivalently characterized as the smallest σ -algebra on T , which contains all closed subsets of T . Measures on $\mathcal{B}(T)$, which are locally finite and inner regular, are commonly referred to as *Radon measures*. To be more specific, we shall call a Radon measure on $\mathcal{B}(T)$ a *T -Radon measure*. A T -Radon measure μ is called a *probability T -Radon measure*, if $\mu(T) = 1$. The set pr_T consisting of all probability T -Radon measures is hence given by

$$\text{pr}_T := \{\mu : \mathcal{B}(T) \rightarrow [0, \infty] \mid \mu \text{ locally finite, inner regular; } \mu(T) = 1\}. \quad (6.3)$$

We also define for every Borel set $S \in \mathcal{B}(T)$ the set $\text{pr}_T(S)$, which consists of all probability T -Radon measures that are supported on S . More precisely, the set $\text{pr}_T(S)$ is given by

$$\text{pr}_T(S) := \{\mu \in \text{pr}_T \mid \mu(S) = 1\}. \quad (6.4)$$

Finally, if T is a sequentially compact topological space, we define for every continuous functional $g : T \rightarrow \mathbb{R}$ the set

$$\arg \min_T g := \left\{ t_* \in T \mid g(t_*) = \min_{t \in T} g(t) \right\}, \quad (6.5)$$

$$\arg \max_T g := \left\{ t_* \in T \mid g(t_*) = \max_{t \in T} g(t) \right\} \quad (6.6)$$

The set defined by (6.5) is called the *set of minimizers of g on T* , and the set defined by (6.6) is called the *set of maximizers of g on T* . Notice that neither of the two sets is empty due to the assumed compactness of T and the assumed continuity of g .

With the above definitions at hand, we are able to state the following theorem, which corresponds to Theorem 2.8.2 in [22].

Theorem 6.13. *Let T be a sequentially compact, separable topological Hausdorff space, and let $g : T \times X \rightarrow \mathbb{R}$ be a functional. Let $x \in X$ be a point, such that the image of $g(\cdot, x)$ is bounded. Furthermore, let U be a neighbourhood of x , and let $c > 0$ be a positive constant, such that for each $y \in U$ the functional $g(\cdot, y)$ is continuous, and such that for each $t \in T$ the functional $g(t, \cdot)$ is c -Lipschitz continuous on U . Then, the functional $f : X \rightarrow \mathbb{R}$, defined by*

$$f := \max_{t \in T} g(t, \cdot)$$

is c -Lipschitz continuous near x , and we have

$$\partial f(x) \subseteq \left\{ \int_T q_t(\cdot) \, d\mu(t) \mid q_t \in \partial g(t, \cdot)(x), \mu \in \text{pr}_T \left(\arg \max_T g(\cdot, x) \right) \right\}.$$

Now, suppose that T is a compact subset of \mathbb{R}^r for some $r \in \mathbb{R}$ endowed with the relative topology of \mathbb{R}^r . Given a point $x \in X$, let us further assume that the functional $g(\cdot, x)$ in Theorem 6.13 attains its maximum on a finite set of points, i.e.,

$$\arg \max_T g(\cdot, x) = \{t_{*1}, \dots, t_{*N}\}$$

for some $N \in \mathbb{N}$. Then, it follows from (6.4) that the set of probability T -Radon measures on T , which are supported on the set of maximizers of $g(\cdot, x)$ on T , is given by

$$\text{pr}_T \left(\arg \max_T g(\cdot, x) \right) = \overline{\text{conv}} \{ \delta_{t_{*1}}, \dots, \delta_{t_{*N}} \}. \quad (6.7)$$

By δ_t we denote the so-called *Dirac measure* on $\mathcal{B}(T)$, which is centered on the point $t \in T$. This measure is defined by

$$\delta_t(S) = \begin{cases} 1 & \text{if } t \in S, \\ 0 & \text{else} \end{cases} \quad \text{for all } S \in \mathcal{B}(T).$$

One easily verifies that δ_t is locally finite and inner regular, and hence a T -Radon measure, for all $t \in T$. The inner regularity follows from the fact that singletons, i.e. sets which contain exactly one point, are compact in \mathbb{R}^r . Furthermore, δ_t is

a probability T -Radon measure for all $t \in T$, since $\delta_t(T) = 1$. Since we also have $\delta_t(\{t\}) = 1$ for all $t \in T$, it follows that every measure in $\text{pr}_T(\{t_{*1}, \dots, t_{*N}\})$ is given by a convex combination of the Dirac measures, which are centered on t_{*1}, \dots, t_{*N} . This is precisely the assertion of (6.7).

It is well-known, that for every continuous functional $h : T \rightarrow \mathbb{R}$ we have

$$\int_T h(t) d\delta_{t_0}(t) = h(t_0) \quad \text{for all } t_0 \in T. \quad (6.8)$$

From (6.8) and the discussion above we deduce the following corollary to Theorem 6.13.

Corollary 6.14. *Let T be a compact subset of \mathbb{R}^r for some $r \in \mathbb{N}$. Under the assumptions of Theorem 6.13, we further assume that*

$$\arg \max_T g(\cdot, x) = \{t_{*1}, \dots, t_{*N}\},$$

for some number $N \in \mathbb{N}$. Then, we have that the generalized differential of the functional f at the point x is given by

$$\partial f(x) \subseteq \left\{ \sum_{l=1}^N \theta_l q_l \mid q_l \in \partial g(t_{*l}, \cdot)(x), \theta_l \in [0, 1], \sum_{l=1}^N \theta_l = 1 \right\}.$$

Another result, which is related to generalized gradients of point-wise maxima, is stated by the following proposition. This proposition corresponds to Proposition 2.3.12 in [22].

Proposition 6.15. *Given a finite number of functionals $g_1, \dots, g_N : X \rightarrow \mathbb{R}$ for some $n \in \mathbb{N}$, let $f : X \rightarrow \mathbb{R}$ be defined by*

$$f(x) := \max\{g_l(x) \mid l = 1, \dots, N\}.$$

If the functionals g_1, \dots, g_N are Lipschitz continuous near a point $x \in X$, then f is Lipschitz continuous near x , and we have

$$\partial f(x) \subseteq \left\{ \sum_{l \in I(x)} \theta_l q_l \mid q_l \in \partial g_l(x), \theta_l \in [0, 1], \sum_{l \in I(x)} \theta_l = 1 \right\},$$

where the index set $I(x)$ is given by

$$I(x) := \{l \in \{1, \dots, N\} \mid g_l(x) = f(x)\}.$$

Finally, we remark that some of the inclusion relations in the propositions and theorems above can be replaced by identities under additional assumptions on the functionals. In Proposition 6.10, for example, equality holds between the two sets if all but at most one of the functionals f_1, \dots, f_N are strictly differentiable at x (see Corollary 1 to Proposition 2.3.3 in [22]). According to Proposition 6.7(d) strict differentiability holds at a point $x \in X$, if and only if the generalized differentials at x reduce to singletons. As we will see in the following section, this condition does not hold in the case of the gap width functionals. For this reason, we chose to present the above propositions and theorems in their most general form.

6.3 Differentiability of the Gap Width Functionals

In Section 6.1 and Section 6.2 we presented a number of general results related to generalized differentials of locally Lipschitz continuous functionals. In Section 5.2 we showed that the goal functionals of photonic band structure optimization problems are typically locally Lipschitz continuous. In Section 5.5 we defined in particular for every index $j \in \mathbb{N}$ the locally Lipschitz continuous gap width functional $J_j : \mathcal{D} \rightarrow \mathbb{R}$ by

$$J_j(\rho) := \max_{\mathbf{k} \in \mathbb{B}} \omega_j(\rho, \mathbf{k}) - \min_{\mathbf{k} \in \mathbb{B}} \omega_{j+1}(\rho, \mathbf{k}).$$

The aim of this section is to characterize the generalized differentials of these functionals.

We begin with some preliminary definitions. Given a coefficient $\rho \in \mathcal{D}$, a vector $\mathbf{k} \in \mathbb{B}$ and an index $j \in \mathbb{N}$ we define the function space

$$\mathbf{U}_j(\rho, \mathbf{k}) := \text{span}\{\mathbf{u}_i(\rho, \mathbf{k}) \mid \lambda_i(\rho, \mathbf{k}) = \lambda_j(\rho, \mathbf{k})\}.$$

Recall that $\lambda_j(\rho, \mathbf{k})$ denotes the j -th smallest eigenvalue of the eigenvalue equation (4.42), and that $\mathbf{u}_j(\rho, \mathbf{k})$ denotes its corresponding function. Hence, the function space $\mathbf{U}_j(\rho, \mathbf{k})$ is the eigenspace corresponding to $\lambda_j(\rho, \mathbf{k})$. By Proposition 4.20, every eigenspace $\mathbf{U}_j(\rho, \mathbf{k})$ is finite-dimensional. Recall also that every eigenfunction $\mathbf{u}_j(\rho, \mathbf{k})$ is an element of the complex Hilbert space $\mathbf{V}_{\mathbf{k}}$, which we defined in Section 4.3. Hence, $\mathbf{U}_j(\rho, \mathbf{k})$ is a finite-dimensional, linear subspace of $\mathbf{V}_{\mathbf{k}}$ for all $\rho \in \mathcal{D}$, $\mathbf{k} \in \mathbb{B}$ and $j \in \mathbb{N}$.

Next, we define for given $\rho \in \mathcal{D}$, $\mathbf{k} \in \mathbb{B}$ and $j \in \mathbb{N}$ the set of functions

$$\mathbf{U}_j^1(\rho, \mathbf{k}) := \{\mathbf{u} \in \mathbf{U}_j(\rho, \mathbf{k}) \mid \|\mathbf{u}\|_{\Omega} = 1\}.$$

Clearly, the set $\mathbf{U}_j^1(\rho, \mathbf{k})$ is given by the intersection of the eigenspace $\mathbf{U}_j(\rho, \mathbf{k})$ with the unit sphere in the complex Hilbert space \mathbf{Z} . The space \mathbf{Z} was defined in Section 4.3.

With the above definitions, we have the following fundamental result.

Theorem 6.16. *For every index $j \in \mathbb{N}$ and every vector $\mathbf{k} \in \mathbb{B}$ the generalized differentials of the functional $\lambda_j(\cdot, \mathbf{k}) : \mathcal{D} \rightarrow \mathbb{R}$ are given by*

$$\partial\lambda_j(\cdot, \mathbf{k})(\rho) = \left\{ p \in \mathcal{E}^* \mid p(\eta) = \int_{\Omega} \eta |(\nabla + i\mathbf{k}) \times \mathbf{u}|^2, \mathbf{u} \in \mathbf{U}_j^1(\rho, \mathbf{k}) \right\}$$

for all $\rho \in \mathcal{D}$.

Proof. According to Theorem 5.2 on page 93 the functional $\lambda_j(\cdot, \mathbf{k}) : \mathcal{D} \rightarrow \mathbb{R}$ is locally Lipschitz continuous for every index $j \in \mathbb{N}$ and every vector $\mathbf{k} \in \mathbb{B}$. Hence, the generalized differential of $\lambda_j(\cdot, \mathbf{k})$ exists at every point $\rho \in \mathcal{D}$.

Now let $\rho_0 \in \mathcal{D}$ be chosen arbitrarily. By definition of the set \mathcal{D} (see (4.32) on page 59) we have that

$$\operatorname{ess\,inf}_{\Omega}(\rho_0) \geq \delta > 0$$

for some positive number $\delta > 0$. Letting $\mathcal{B}_{\delta/2}(\rho_0)$ denote the open ball in \mathcal{E} with radius $\delta/2$ and center ρ_0 , we have that $\mathcal{B}_{\delta/2}(\rho_0) \subset \mathcal{D}$, as well as

$$\operatorname{ess\,sup}_{\Omega}(\rho_0) + \delta \geq \operatorname{ess\,sup}_{\Omega}(\rho) \geq \operatorname{ess\,inf}_{\Omega}(\rho) \quad \text{for all } \rho \in \mathcal{B}_{\delta/2}(\rho_0).$$

Next, we define the function $\widehat{\kappa} : \mathbb{B} \rightarrow \mathbb{R}_{\geq 0}$ by

$$\widehat{\kappa}(\mathbf{k}) := \left(\operatorname{ess\,sup}_{\Omega}(\rho_0) + \delta \right) \kappa_0(\mathbf{k}) \quad \text{for all } \mathbf{k} \in \mathbb{B},$$

where κ_0 denotes the function defined in Lemma 4.4 on page 52. By definition of the function κ (see (4.40) on page 63) and the above estimate, we have that $\widehat{\kappa}(\mathbf{k}) \geq \kappa(\rho, \mathbf{k})$ for all $\rho \in \mathcal{B}_{\delta/2}(\rho_0)$ and all $\mathbf{k} \in \mathbb{B}$.

For every $\rho \in B_{\delta/2}^{\infty}(\rho_0)$ and every vector $\mathbf{k} \in \mathbb{B}$ we define the sesquilinear form $\widetilde{l}_{\mathbf{k}}(\rho) : \mathbf{V}_{\mathbf{k}} \times \mathbf{V}_{\mathbf{k}} \rightarrow \mathbb{C}$ by $\widetilde{l}_{\mathbf{k}}(\rho) := a_{\mathbf{k}}(\rho) + \widehat{\kappa}(\mathbf{k})m$. Recall that the sesquilinear forms $a_{\mathbf{k}}(\rho)$ and m were defined by (4.33) and (4.35) on page 59. According to Proposition 4.14, for every $\rho \in B_{\delta/2}^{\infty}(\rho_0)$ and every $\mathbf{k} \in \mathbb{B}$ we have that

$$\widetilde{l}_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{u}) \geq a_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{u}) + \kappa(\rho, \mathbf{k})m(\mathbf{u}, \mathbf{u}) \geq \alpha(\rho, \mathbf{k})\|\mathbf{u}\|_{\operatorname{curl}, \Omega}^2 \quad \text{for all } \mathbf{u} \in \mathbf{V}_{\mathbf{k}}.$$

It follows that the sesquilinear form $\widetilde{l}_{\mathbf{k}}(\rho)$ is conjugate-symmetric, bounded and $\mathbf{V}_{\mathbf{k}}$ -elliptic for all $\rho \in B_{\delta/2}^{\infty}(\rho_0)$ and all $\mathbf{k} \in \mathbb{B}$.

Given a vector $\mathbf{k} \in \mathbb{B}$, we denote by $\overline{\mathbb{L}}(\mathbf{V}_{\mathbf{k}})$ the linear space consisting of all continuous conjugate linear operators from $\mathbf{V}_{\mathbf{k}}$ into $\mathbf{V}_{\mathbf{k}}^*$. Furthermore, let $\mathbb{E}(\mathbf{V}_{\mathbf{k}})$ be the set of all operators $T \in \overline{\mathbb{L}}(\mathbf{V}_{\mathbf{k}})$, such that the mapping $(\mathbf{u}, \mathbf{v}) \mapsto T\mathbf{u}[\mathbf{v}]$ is

a conjugate-symmetric, bounded, $\mathbf{V}_{\mathbf{k}}$ -elliptic sesquilinear form on $\mathbf{V}_{\mathbf{k}} \times \mathbf{V}_{\mathbf{k}}$. With this, we define for every $\mathbf{k} \in \mathbb{B}$ the operator-valued function $\tilde{L}_{\mathbf{k}} : \mathcal{B}_{\delta/2}(\rho_0) \rightarrow \mathbb{E}(\mathbf{V}_{\mathbf{k}})$ by

$$\tilde{L}_{\mathbf{k}}(\rho)\mathbf{u}[\mathbf{v}] := \tilde{l}_{\mathbf{k}}(\rho)(\mathbf{u}, \mathbf{v}) \quad \text{for all } \rho \in \mathcal{B}_{\delta/2}(\rho_0), \mathbf{u}, \mathbf{v} \in \mathbf{V}_{\mathbf{k}}.$$

According to the definition in Section 4.4, we have that $\tilde{L}_{\mathbf{k}}(\rho) = A_{\mathbf{k}}(\rho) + \tilde{\kappa}(\mathbf{k})M_{\mathbf{k}}|_{\mathbf{V}_{\mathbf{k}}}$ for all $\rho \in \mathcal{B}_{\delta/2}(\rho_0)$ and all $\mathbf{k} \in \mathbb{B}$. It follows that the j -th smallest eigenvalue $\tilde{\lambda}_j(\rho, \mathbf{k})$ of $\tilde{L}_{\mathbf{k}}(\rho)$ coincides with $\lambda_j(\rho, \mathbf{k}) + \tilde{\kappa}(\mathbf{k})$ for all $\rho \in \mathcal{B}_{\delta/2}(\rho_0)$ and all $\mathbf{k} \in \mathbb{B}$. Furthermore, we have that the eigenspaces of $\tilde{\lambda}_j(\rho, \mathbf{k})$ and $\lambda_j(\rho, \mathbf{k}) + \tilde{\kappa}(\mathbf{k})$ coincide. Since the function $\tilde{\kappa}$ only depends on \mathbf{k} , we also have that $\partial\tilde{\lambda}_j(\cdot, \mathbf{k})(\rho) = \partial\lambda_j(\cdot, \mathbf{k})(\rho)$ for all $\rho \in \mathcal{B}_{\delta/2}(\rho_0)$ and all $\mathbf{k} \in \mathbb{B}$.

Now let $\mathbf{k} \in \mathbb{B}$ be fixed. Note that for every coefficient $\rho \in \mathcal{B}_{\delta/2}(\rho_0)$ the operator $\tilde{L}_{\mathbf{k}}(\rho)$ satisfies

$$\tilde{L}_{\mathbf{k}}(\rho)\mathbf{u}[\mathbf{v}] = \int_{\Omega} \overline{\rho(\nabla + i\mathbf{k}) \times \mathbf{u}} \cdot (\nabla + i\mathbf{k}) \times \mathbf{v} + \tilde{\kappa}(\mathbf{k}) \int_{\Omega} \overline{\mathbf{u}} \cdot \mathbf{v} \quad \text{for all } \mathbf{u}, \mathbf{v} \in \mathbf{V}_{\mathbf{k}}.$$

Given a function $\eta \in \mathcal{E}$, such that $\rho_0 + \eta \in \mathcal{B}_{\delta/2}(\rho_0)$, we thus find that

$$\begin{aligned} (\tilde{L}_{\mathbf{k}}(\rho_0 + \eta) - \tilde{L}_{\mathbf{k}}(\rho_0))\mathbf{u}[\mathbf{v}] &= \int_{\Omega} \eta \overline{(\nabla + i\mathbf{k}) \times \mathbf{u}} \cdot (\nabla + i\mathbf{k}) \times \mathbf{v} \\ &= A_{\mathbf{k}}(\eta)\mathbf{u}[\mathbf{v}] \quad \text{for all } \mathbf{u}, \mathbf{v} \in \mathbf{V}_{\mathbf{k}}. \end{aligned}$$

Clearly, the mapping $\eta \mapsto A_{\mathbf{k}}(\eta)$ is linear from \mathcal{E} into $\overline{\mathbb{L}}(\mathbf{V}_{\mathbf{k}})$. From this we deduce that the function $\tilde{L}_{\mathbf{k}}$ is continuously Gâteaux differentiable at ρ_0 with Gâteaux derivative $D\tilde{L}_{\mathbf{k}}(\rho_0) = A_{\mathbf{k}}(\cdot)$.

We have thus shown that the operator-valued function $\tilde{L}_{\mathbf{k}}$ satisfies the hypotheses in Section 1 in [28]. The assertion of this theorem then follows directly from Theorem 1 in [28]. \square

According to (4.74) on page 80 the functions $\omega_j : \mathcal{D} \times \mathbb{B} \rightarrow \mathbb{R}$ defined as $\omega_j = \sqrt{\lambda_j}$ for all $j \in \mathbb{N}$. By Theorem 6.11 we have that the generalized differentials of these functionals are given by

$$\partial\omega_j(\cdot, \mathbf{k})(\rho) = \left\{ p \in \mathcal{E}^* \left| p(\eta) = \int_{\Omega} \eta \frac{|(\nabla + i\mathbf{k}) \times \mathbf{u}|^2}{2\omega_j(\rho, \mathbf{k})}, \right. \right. \\ \left. \left. \mathbf{u} \in \mathbf{U}_j^1(\rho, \mathbf{k}) \right\} \quad (6.9)$$

for all $j \in \mathbb{N}$, $\rho \in \mathcal{D}$ and $\mathbf{k} \in \mathbb{B}$. For convenience, we define for every index $j \in \mathbb{N}$ the functionals $J_{\text{up},j}, J_{\text{lo},j} : \mathcal{D} \rightarrow \mathbb{R}$ by

$$J_{\text{lo},j}(\rho) := \max_{\mathbf{k} \in \mathbb{B}} \omega_j(\rho, \mathbf{k}), \quad \text{for all } j \in \mathbb{N}, \rho \in \mathcal{D}, \quad (6.10)$$

$$J_{\text{up},j}(\rho) := \min_{\mathbf{k} \in \mathbb{B}} \omega_{j+1}(\rho, \mathbf{k}), \quad \text{for all } j \in \mathbb{N}, \rho \in \mathcal{D}. \quad (6.11)$$

Given an index j , we call the functional $J_{\text{up},j}$ the j -th upper band edge functional. Similarly, we refer to $J_{\text{lo},j}$ as the j -th lower band edge functional. By Theorem 6.13 we obtain

$$\partial J_{\text{lo},j}(\rho) \subseteq \left\{ \int_{\mathbb{B}} q_{\mathbf{k}}(\cdot) d\mu(\mathbf{k}) \left| \begin{array}{l} q_{\mathbf{k}} \in \partial \omega_j(\cdot, \mathbf{k})(\rho), \\ \mu \in \text{pr}_{\mathbb{B}} \left(\arg \max_{\mathbb{B}} \omega_j(\rho, \cdot) \right) \end{array} \right. \right\}, \quad (6.12)$$

$$\partial J_{\text{up},j}(\rho) \subseteq \left\{ \int_{\mathbb{B}} q_{\mathbf{k}}(\cdot) d\mu(\mathbf{k}) \left| \begin{array}{l} q_{\mathbf{k}} \in \partial \omega_{j+1}(\cdot, \mathbf{k})(\rho), \\ \mu \in \text{pr}_{\mathbb{B}} \left(\arg \min_{\mathbb{B}} \omega_{j+1}(\rho, \cdot) \right) \end{array} \right. \right\} \quad (6.13)$$

for all $j \in \mathbb{N}$ and all $\rho \in \mathcal{D}$. Since $J_j = J_{\text{lo},j} - J_{\text{up},j}$, we have by Proposition 6.9 and Proposition 6.10 that the generalized differential of the j -th gap width functional satisfies the inclusion relation

$$\partial J_j(\rho) \subseteq \partial J_{\text{lo},j}(\rho) - \partial J_{\text{up},j}(\rho) \quad \text{for all } j \in \mathbb{N}, \rho \in \mathcal{D}. \quad (6.14)$$

In practice, one finds that the functionals $\omega_j(\rho, \cdot)$ typically attain their extrema at a small number of points in the first Brillouin zone \mathbb{B} . Hence, by assuming that the set of minimizers and maximizers are given by

$$\arg \min_{\mathbb{B}} \omega_j(\rho, \cdot) = \{ \mathbf{k}_{\text{lo},j}^{(1)}, \dots, \mathbf{k}_{\text{lo},j}^{(N_{\text{lo},j})} \}, \quad (6.15)$$

$$\arg \max_{\mathbb{B}} \omega_{j+1}(\rho, \cdot) = \{ \mathbf{k}_{\text{up},j}^{(1)}, \dots, \mathbf{k}_{\text{up},j}^{(N_{\text{up},j})} \} \quad (6.16)$$

for some numbers $N_{\text{lo},j}, N_{\text{up},j} \in \mathbb{N}$, we obtain by Corollary 6.14 the inclusion relations

$$\partial J_{\text{lo},j}(\rho) \subseteq \left\{ \sum_{l=1}^{N_{\text{lo},j}} \theta_l q_l \left| \begin{array}{l} q_l \in \partial \omega_j(\cdot, \mathbf{k}_{\text{lo},j}^{(l)})(\rho), \theta_l \in [0, 1], \\ \sum_{l=1}^{N_{\text{lo},j}} \theta_l = 1 \end{array} \right. \right\}, \quad (6.17)$$

$$\partial J_{\text{up},j}(\rho) \subseteq \left\{ \sum_{l=1}^{N_{\text{up},j}} \theta_l q_l \left| \begin{array}{l} q_l \in \partial \omega_{j+1}(\cdot, \mathbf{k}_{\text{up},j}^{(l)})(\rho), \theta_l \in [0, 1], \\ \sum_{l=1}^{N_{\text{up},j}} \theta_l = 1 \end{array} \right. \right\} \quad (6.18)$$

for all $j \in \mathbb{N}$ and all $\rho \in \mathcal{D}$.

For our further discussion, we find it useful to define for every function $\mathbf{u} \in \mathbf{V}_{\mathbf{k}}$, every index $j \in \mathbf{V}_{\mathbf{k}}$ the functions $\gamma_{\text{lo},j}^{(l)}(\mathbf{u}), \gamma_{\text{up},j}^{(l)}(\mathbf{u}) : \Omega \rightarrow \mathbb{R}$ by

$$\gamma_{\text{lo},j}^{(l)}(\mathbf{u}) := \frac{|(\nabla + i\mathbf{k}_{\text{lo},j}^{(l)}) \times \mathbf{u}|^2}{2 \omega_j(\rho, \mathbf{k}_{\text{lo},j}^{(l)})} \quad \text{for } l = 1, \dots, N_{\text{lo},j},$$

$$\gamma_{\text{up},j}^{(l)}(\mathbf{u}) := \frac{|(\nabla + \mathbf{i}\mathbf{k}_{\text{up},j}^{(l)}) \times \mathbf{u}|^2}{2\omega_{j+1}(\rho, \mathbf{k}_{\text{up},j}^{(l)})} \quad \text{for } l = 1, \dots, N_{\text{up},j}.$$

Notice that $(\nabla + \mathbf{i}\mathbf{k}) \times \mathbf{u}$ is a vector field in $\mathbf{Z} = L^2(\Omega)^3$ for every vector $\mathbf{k} \in \mathbb{B}$ and every vector field $\mathbf{u} \in \mathbf{V}_{\mathbf{k}}$ according to Proposition 4.3(a) on page 50. It thus follows, that $|(\nabla + \mathbf{i}\mathbf{k}) \times \mathbf{u}|^2$ is a scalar field in $L^1(\Omega, \mathbb{R})$ for all $\mathbf{k} \in \mathbb{B}$ and all $\mathbf{u} \in \mathbf{V}_{\mathbf{k}}$, and the same also applies to the functions $\gamma_{\text{lo},j}^{(l)}(\mathbf{u})$ and $\gamma_{\text{up},j}^{(l)}(\mathbf{u})$.

According to Theorem 6.16 and (6.9), the inclusion relations (6.17) and (6.18) can be written equivalently as

$$\partial J_{\text{lo},j}(\rho) \subseteq \left\{ p \in \mathcal{E}^* \mid p(\eta) = \int_{\Omega} \eta \gamma, \gamma \in \Gamma_{\text{lo},j}(\rho) \right\}, \quad (6.19)$$

$$\partial J_{\text{up},j}(\rho) \subseteq \left\{ p \in \mathcal{E}^* \mid p(\eta) = \int_{\Omega} \eta \gamma, \gamma \in \Gamma_{\text{up},j}(\rho) \right\} \quad (6.20)$$

for all $j \in \mathbb{N}$ and all $\rho \in \mathcal{D}$, where the sets $\Gamma_{\text{lo},j}(\rho)$ and $\Gamma_{\text{up},j}(\rho)$ are defined by

$$\Gamma_{\text{lo},j}(\rho) := \left\{ \sum_{l=1}^{N_{\text{lo},j}} \theta_l \gamma_{\text{lo},j}^{(l)}(\mathbf{u}^{(l)}) \mid \mathbf{u}^{(l)} \in \mathbf{U}_j^1(\rho, \mathbf{k}_{\text{lo},j}^{(l)}), \right. \\ \left. \theta_l \in [0, 1], \sum_{l=1}^{N_{\text{lo},j}} \theta_l = 1 \right\},$$

$$\Gamma_{\text{up},j}(\rho) := \left\{ \sum_{l=1}^{N_{\text{up},j}} \theta_l \gamma_{\text{up},j}^{(l)}(\mathbf{u}^{(l)}) \mid \mathbf{u}^{(l)} \in \mathbf{U}_{j+1}^1(\rho, \mathbf{k}_{\text{up},j}^{(l)}), \right. \\ \left. \theta_l \in [0, 1], \sum_{l=1}^{N_{\text{up},j}} \theta_l = 1 \right\}.$$

Under the hypotheses (6.15) and (6.16), the inclusion relations (6.19) and (6.20) reveal some important aspects about the generalized differentials of the gap width functionals. First, we notice that $\Gamma_{\text{lo},j}(\rho)$ and $\Gamma_{\text{up},j}(\rho)$ are subsets of $L^1(\Omega, \mathbb{R})$. Hence, according to (6.14), every continuous, linear functional contained in $\partial J_j(\rho)$ can be expressed as a weighted integral over Ω with a weight function in $L^1(\Omega, \mathbb{R})$. The generalized differentials of the gap width functionals are hence weak*-closed, convex, bounded subsets of

$$\mathcal{E}_{\text{reg}}^* := \left\{ p \in \mathcal{E}^* \mid p(\eta) = \int_{\Omega} \eta w, w \in L^1(\Omega, \mathbb{R}) \right\}.$$

Since the Banach space \mathcal{E} is not reflexive, the space $\mathcal{E}_{\text{reg}}^*$ is a proper linear subspace of \mathcal{E}^* . It is easy to see that the space $\mathcal{E}_{\text{reg}}^*$ is isometrically, linearly isomorphic to $L^1(\Omega, \mathbb{R})$.

Concerning the differentiability of gap width functionals, we are able to make the following conclusive statement. Recall that by Proposition 6.7(d) a gap width functional J_j is strictly differentiable at a point $\rho \in \mathcal{D}$, if and only if the generalized differential $\partial J_j(\rho)$ contains exactly one element, which is then the strict differential of J_j at ρ . According to (6.19) and (6.20), the generalized differential $\partial J_j(\rho)$ is given by a singleton, if and only if each of the functions $\omega_j(\rho, \cdot)$ and $\omega_{j+1}(\rho, \cdot)$ attain their minimum and maximum at exactly one vector $\mathbf{k}_{\text{lo},j} \in \mathbb{B}$ and $\mathbf{k}_{\text{up},j} \in \mathbb{B}$, respectively, and if the corresponding eigenspaces $\mathbf{U}_j(\rho, \mathbf{k}_{\text{lo},j})$ and $\mathbf{U}_{j+1}(\rho, \mathbf{k}_{\text{lo},j})$ are one-dimensional.

Finally, we remark that the sets $\Gamma_{\text{lo},j}(\rho)$ and $\Gamma_{\text{up},j}(\rho)$ are completely determined by eigensolutions of the weakly formulated family of eigenvalue problems given by Problem 4.16. These eigensolutions also need to be computed in order to evaluate the gap width functional J_j at a point $\rho \in \mathcal{D}$. Therefore, in a numerical algorithm the sets $\Gamma_{\text{lo},j}(\rho)$ and $\Gamma_{\text{up},j}(\rho)$ can be computed simultaneously with the goal value $J_j(\rho)$.

6.4 The Two-Dimensional Case

In the previous section we derived the generalized differentials of the gap width functionals J_j for three-dimensional photonic crystals. In Section 5.5 we also defined for every index $j \in \mathbb{N}$ the TM and TE gap width functionals $J^{\text{TE}} : \mathcal{D} \rightarrow \mathbb{R}$ and $J^{\text{TE}} : \mathcal{D} \rightarrow \mathbb{R}$ by

$$\begin{aligned} J_j^{\text{TM}}(\rho) &:= \max_{\mathbf{k} \in \mathbb{B}} \omega_j^{\text{TM}}(\rho, \mathbf{k}) - \min_{\mathbf{k} \in \mathbb{B}} \omega_{j+1}^{\text{TM}}(\rho, \mathbf{k}) && \text{for all } \rho \in \mathcal{D}, \\ J_j^{\text{TE}}(\rho) &:= \max_{\mathbf{k} \in \mathbb{B}} \omega_j^{\text{TE}}(\rho, \mathbf{k}) - \min_{\mathbf{k} \in \mathbb{B}} \omega_{j+1}^{\text{TE}}(\rho, \mathbf{k}) && \text{for all } \rho \in \mathcal{D}. \end{aligned}$$

For completeness, we give a characterization of the generalized differentials of these functionals in this section.

In the following we assume that Ω is a primitive domain of a Bravais lattice in \mathbb{R}^2 of rank 2. By \mathbb{B} we denote the lattice's first Brillouin zone. In analogy to the three-dimensional setting, we define the eigenspaces

$$\begin{aligned} U_j^{\text{TM}}(\rho, \mathbf{k}) &:= \text{span}\{u_i^{\text{TM}}(\rho, \mathbf{k}) \mid \lambda_i(\rho, \mathbf{k}) = \lambda_j(\rho, \mathbf{k})\}, \\ U_j^{\text{TE}}(\rho, \mathbf{k}) &:= \text{span}\{u_i^{\text{TE}}(\rho, \mathbf{k}) \mid \lambda_i^{\text{TE}}(\rho, \mathbf{k}) = \lambda_j^{\text{TE}}(\rho, \mathbf{k})\}, \end{aligned}$$

as well as the function sets

$$\begin{aligned} U_j^{\text{TM}^1}(\rho, \mathbf{k}) &:= \{u \in U_j^{\text{TM}}(\rho, \mathbf{k}) \mid \|u\|_\rho = 1\}, \\ U_j^{\text{TE}^1}(\rho, \mathbf{k}) &:= \{u \in U_j^{\text{TE}}(\rho, \mathbf{k}) \mid \|u\|_\Omega = 1\}. \end{aligned}$$

for all $j \in \mathbb{N}$, $\rho \in \mathcal{D}$, $\rho \in \mathcal{D}$ and all $\mathbf{k} \in \mathbb{B}$. Recall that the norm $\|\cdot\|_\rho$ was defined by $\|u\|_\rho^2 = m^{\text{TM}}(\rho)(u, u)$ for all $u \in L^2(\Omega)$ (see (4.87) on page 86). With the above definition, we can give the following characterization of the generalized differentials of the TM and TE eigenvalues.

Theorem 6.17. *For every index $j \in \mathbb{N}$ and every vector $\mathbf{k} \in \mathbb{B}$ the generalized differentials of the functionals $\lambda_j^{\text{TM}}(\cdot, \mathbf{k}) : \mathcal{D} \rightarrow \mathbb{R}$ and $\lambda_j^{\text{TE}}(\cdot, \mathbf{k}) : \mathcal{D} \rightarrow \mathbb{R}$ are given by*

$$\begin{aligned} \partial\lambda_j^{\text{TM}}(\cdot, \mathbf{k})(\rho) &= \left\{ p \in \mathcal{E}^* \left| p(\eta) = -\lambda_j^{\text{TM}}(\rho, \mathbf{k}) \int_{\Omega} \eta |u|^2, u \in U_j^{\text{TM}1}(\rho, \mathbf{k}) \right. \right\}, \\ \partial\lambda_j^{\text{TE}}(\cdot, \mathbf{k})(\rho) &= \left\{ p \in \mathcal{E}^* \left| p(\eta) = \int_{\Omega} \eta |(\nabla + i\mathbf{k})u|^2, u \in U_j^{\text{TE}1}(\rho, \mathbf{k}) \right. \right\} \end{aligned}$$

for all $\rho \in \mathcal{D}$ and all $\rho \in \mathcal{D}$.

A proof of the first identity in Theorem 6.17 is given in Section 3 in [29]. The second identity can be proved by a similar line of argument as given in the proof of Theorem 6.16. According to Theorem 6.11, the generalized differentials of the functions $\omega_j^{\text{TM}} : \mathcal{D} \rightarrow \mathbb{R}$ and $\omega_j^{\text{TE}} : \mathcal{D} \rightarrow \mathbb{R}$ are given by

$$\begin{aligned} \partial\omega_j^{\text{TM}}(\cdot, \mathbf{k})(\rho) &= \left\{ p \in \mathcal{E}^* \left| p(\eta) = -\omega_j^{\text{TM}}(\rho, \mathbf{k}) \int_{\Omega} \eta |u|^2, \right. \right. \\ &\quad \left. \left. u \in U_j^{\text{TM}1}(\rho, \mathbf{k}) \right. \right\}, \\ \partial\omega_j^{\text{TE}}(\cdot, \mathbf{k})(\rho) &= \left\{ p \in \mathcal{E}^* \left| p(\eta) = \int_{\Omega} \eta \frac{|(\nabla + i\mathbf{k})u|^2}{2\omega_j^{\text{TE}}(\rho, \mathbf{k})}, \right. \right. \\ &\quad \left. \left. u \in U_j^{\text{TE}1}(\rho, \mathbf{k}) \right. \right\} \end{aligned}$$

for all $j \in \mathbb{N}$, $\rho \in \mathcal{D}$, and all $\mathbf{k} \in \mathbb{B}$. From here, one can proceed exactly as in the previous section and characterize the generalized differentials of J_j^{TM} and J_j^{TE} in terms of integrals with respect to the variable \mathbf{k} over subsets of \mathbb{B} , where the functions ω_j^{TM} and ω_j^{TE} attain their minima or maxima. By defining the *TM band edge functionals* $J_{\text{lo},j}^{\text{TM}}, J_{\text{up},j}^{\text{TM}} : \mathcal{D} \rightarrow \mathbb{R}$, as well as the *TE band edge functionals* $J_{\text{lo},j}^{\text{TE}}, J_{\text{up},j}^{\text{TE}} : \mathcal{D} \rightarrow \mathbb{R}$ in analogy to (6.10) and (6.11), one obtains inclusion relations for their generalized differentials, which are analogous to (6.12) and (6.13). By assuming that the minima and maxima are attained at a finite number of quasimomentum vectors \mathbf{k} , one obtains inclusion relations analogous to (6.17) and (6.18). The generalized gradients of the TM gap width functionals J_j^{TM} satisfy the inclusion relations $\partial J_j^{\text{TM}}(\rho) \subseteq \partial J_{\text{lo},j}^{\text{TM}}(\rho) - J_{\text{up},j}^{\text{TM}}(\rho)$ for all $j \in \mathbb{N}$ and all $\rho \in \mathcal{D}$, and the analogous result also holds for the generalized gradients of the TE gap width functionals J_j^{TE} .

Chapter 7

A Generalized Gradient Method

In this section, we develop an optimization algorithm for discretized photonic band structure optimization problems based on the concept of generalized gradient methods. In Section 7.1 we outline the basic idea of generalized gradient methods and comment on their mathematical justification and practical applicability. In Section 7.2 we present a discretization scheme for photonic band structure optimization problems. In Section 7.3 and Section 7.4 we develop two approaches to incorporate certain optimization constraints into a generalized gradient method for photonic band structure optimization problems. An iterative optimization algorithm based on these approaches is presented in Section 7.5. A crucial step in each iteration of this algorithm is the choice of a so-called descent direction. In Section 7.6 we present a strategy for the choice of descent directions, which can be applied to photonic band gap maximization problems. In Section 7.7 we make some final remarks on the generalized gradient method.

7.1 Basic Concepts

We consider a minimization problem of the form

$$\underset{x \in X}{\text{minimize}} \quad f(x), \tag{7.1}$$

where f is a locally Lipschitz continuous goal functional defined on a finite-dimensional Hilbert space X over \mathbb{R} with inner product $\langle \cdot, \cdot \rangle_X$. A minimization problem of this type is commonly referred to as a *locally Lipschitz minimization problem*. Due to its local Lipschitz continuity, the functional f admits a generalized differential $\partial f(x)$ at every point $x \in X$. Since X is a Hilbert space, the generalized differential $\partial f(x)$ can be identified with a subset $G_f(x)$ of X via the Rieszian isomorphism. This subset is called the generalized gradient of f at x . It follows from Proposition 6.4 on page 113 that $G_f(x)$ is a closed, convex and bounded

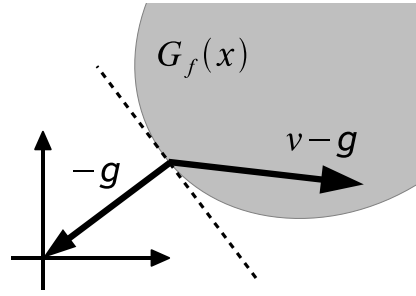


Figure 7.1: The geometric meaning of the variational inequality $\langle -g, v - g \rangle \leq 0$, where v is an arbitrary vector in $G_f(x)$, and where g is the best approximation in $G_f(x)$ of the zero element in X .

subset of X for every $x \in X$. The closedness of $G_f(x)$ follows from the fact that $\partial f(x)$ is weak-*compact, which implies that $G_f(x)$ is weakly compact. Since X is a finite-dimensional Hilbert space, weak compactness implies compactness, which in turn implies closedness according to the Heine–Borel Theorem.

A *generalized gradient method* is an optimization algorithm, which can be used to find local minima of locally Lipschitz continuous functionals defined on a finite-dimensional, real Hilbert space. A generalized gradient method for the minimization problem (7.1) attempts to construct a minimizing sequence x_0, x_1, x_2, \dots , etc. in X for the goal functional f . The *initial point* $x_0 \in X$ is chosen arbitrarily. The points x_1, x_2, \dots etc., are determined iteratively according to

$$x_l := x_{l-1} - s_l g_{l-1} \quad \text{for } l = 1, 2, \dots,$$

where $s_l > 0$ is a suitable *step size*, and where g_{l-1} is a suitable element of $G_f(x_{l-1})$ for every $l \in \mathbb{N}$. Here, the variable $l \in \mathbb{N}_0$ enumerates the iterations of the generalized gradient method. In each iteration, the generalized gradient element g_{l-1} is chosen such that $-g_{l-1}$ is a *descent direction* of f at the point x_{l-1} , if possible. More precisely, $g_{l-1} \in G_f(x_{l-1})$ is chosen such that there exists a positive number $t_{\max} > 0$, such that

$$f(x_{l-1} - t g_{l-1}) \leq f(x_{l-1}) \quad \text{for all } t \in [0, t_{\max}).$$

The following lemma informs us on how an element $g_{l-1} \in G_f(x)$ can be chosen, such that $-g_{l-1}$ is a descent direction.

Lemma 7.1. *Suppose that $0 \notin G_f(x)$. Let $g \in G_f(x)$, such that $\|g\|_X \leq \|v\|_X$ for all $v \in G_f(x)$. Then, $-g$ is a descent direction of f at x .*

Proof. The following proof is a particularization of the proof of Proposition 6.2.4 in [22]. Since $G_f(x)$ is convex and closed, there exists a uniquely defined element $g \in G_f(x)$, such that $\|g\|_X \leq \|v\|_X$ for all $v \in G_f(x)$. Note that the element g is the best approximation in the convex set $G_f(x)$ of the zero element in X . Hence, we have the variational inequality

$$\langle -g, v - g \rangle \leq 0 \quad \text{for all } v \in G_f(x).$$

It should be noted that this variational inequality has a simple geometric meaning, namely that the angle between the vector $-g$ and every other vector, that connects the point g to another point in the set $G_f(x)$, is greater than or equal to $\pi/2$ (see also Figure 7.1).

$$\langle -g, v \rangle_X \leq \langle -g, g \rangle_X = -\|g\|_X^2.$$

By Proposition 2.1.5 in [22], there exists for every given number $\varepsilon > 0$ a number $\delta > 0$, such that

$$G_f(y) \subset G_f(x) + B_\varepsilon(0) \quad \text{for all } y \in X \text{ with } \|x - y\|_X < \delta,$$

where $B_\varepsilon(0)$ denotes the open ball in X with radius ε centered at zero. Here, we choose δ , such that the above inclusion relation holds for $\varepsilon := \|g\|_X/2$. We also set $t_{\max} := \delta/\|g\|_X$. Choosing an arbitrary number $t \in [0, t_{\max})$, we have by Lebourg's theorem (see Theorem 6.6 on page 114) that there exists an element $z \in \text{conv}\{x, x - \delta g\}$ and an element $u \in G_f(z)$, such that

$$f(x - tg) - f(x) = -t\langle u, g \rangle.$$

Since $\|x - z\|_X \leq t\|g\|_X < \delta$, we have that $G_f(z) \subset G_f(x) + B_\varepsilon(0)$. Hence, there exists an element $v \in G_f(z)$ and an element $w \in X$ with $\|w\|_X < \delta$, such that $u = v + w$. Therefore, we have that

$$\begin{aligned} f(x - tg) - f(x) &= -t\langle v + w, g \rangle_X = t\langle -g, v \rangle_X - t\langle w, g \rangle_X \\ &\leq -t\|g\|_X^2 + t\|w\|_X\|g\|_X \leq -t\|g\|_X^2 + t\varepsilon\|g\|_X \\ &= -t\frac{\|g\|_X^2}{2}, \end{aligned}$$

which completes this proof. □

Notice that one needs to solve a convex minimization problem in order to find the element g proposed by Lemma 7.1. If the element g turns out to be zero, then a

necessary condition for a local extremum of f is satisfied at the point x according to Lemma 6.5 on page 113.

If the functional f is convex in a neighbourhood U of the point x , then $-g$ is a descent direction for every element $g \in G_f(x) \setminus \{0\}$. This is due to Proposition 6.7(a) on page 114, which states that $\partial f(x)$ coincides with the subdifferential of f at x under the above assumption. Hence, $G_f(x)$ corresponds to the so-called *subgradient* of f at x . This subgradient is defined as the set of all elements $u \in X$, such that

$$f(y) - f(x) \geq \langle u, y - x \rangle_X \quad \text{for all } y \in U.$$

Finally, recall that $G_f(x)$ is a singleton, if and only if f is strictly differentiable at x . Then, the single element contained in $G_f(x)$ is the strict gradient $\nabla_s f(x)$ of f at x . If the strict gradient does not vanish, it easily follows that $-\nabla_s f(x)$ is a descent direction of f at the point x .

It should be noted, that in many cases the generalized gradient $G_f(x)$ is not known explicitly. Often, one can only determine a superset of $G_f(x)$ using the rules of calculus presented in Section 6.2. In these cases, specific *choice strategies* for the descent directions have to be devised. These choice strategies have a significant impact on the performance of a generalized gradient method. We shall comment on this in Section 7.7.

So far, results concerning the convergence of generalized gradient methods could only be established for specific types of locally Lipschitz minimization problems (see e.g. [79], [19]). As far as can be judged from the literature, a general convergence theory for generalized gradient methods is not available yet.

7.2 Discretization

The proof of Lemma 7.1 relies on the fact that X is a finite-dimensional Hilbert space. In particular, the proof invokes Proposition 2.1.5 in [22], which also assumes that X is finite-dimensional. So far, it is not known whether the results of Lemma 7.1 also apply to infinite-dimensional Hilbert spaces. If X is a Banach space, the additional problem arises that the generalized differential $\partial f(x) \subset X^*$ cannot be identified with a set $G_f(x) \subset X$, in general. In particular, such an identification is generally not possible for $X = \mathcal{E}$. Recall, for example, that the generalized differentials of the gap width functionals can only be identified with subsets of $L^1(\Omega, \mathbb{R})$ (see Section 6.3). Due to this, we were not able to develop a generalized gradient method for the minimization problem

$$\underset{\rho \in \mathcal{C}}{\text{minimize}} \quad J(\rho).$$

However, we could develop a optimization algorithm that can be applied to the corresponding discretized problem

$$\underset{\rho_h \in \mathcal{C}_h}{\text{minimize}} J_h(\rho_h), \quad (7.2)$$

where \mathcal{C}_h denotes a discretization of the admissible set \mathcal{C} , and where J_h denotes the corresponding discretization of the gap width functional J . In this section we explain the underlying discretization scheme

Given a primitive domain Ω of a Bravais lattice $\Lambda \subset \mathbb{R}^3$ of rank 3, we suppose that

$$\mathcal{T}_h = \{T_1, T_2, \dots, T_N\}$$

is a *finite element mesh* on Ω for some $N \in \mathbb{N}$, where the $h > 0$ indicates the maximal diameter of the elements T_i in the mesh. We suppose that the Hilbert spaces \mathbf{W} and Q , as defined by (4.28) and (4.29) on page 59, are discretized by conforming finite element spaces \mathbf{W}_h and Q_h . Given a coefficient $\rho \in \mathcal{D}$ and a finite set of quasimomentum vectors $K_h \subset \mathbb{B}$, an approximate band structure is computed by solving for each $\mathbf{k} \in K_h$ the constrained matrix eigenvalue problem

$$\begin{cases} \mathbf{A}_{\mathbf{k},h}(\rho)\mathbf{u} = \lambda_h \mathbf{M}_{\mathbf{k},h}\mathbf{u}, \\ \mathbf{B}_{\mathbf{k},h}\mathbf{u} = \mathbf{0}. \end{cases} \quad (7.3)$$

Here, $\mathbf{A}_{\mathbf{k},h}(\rho)$, $\mathbf{M}_{\mathbf{k},h}$, and $\mathbf{B}_{\mathbf{k},h}$ are the finite element matrices that realize the operators $A_{\mathbf{k}}(\rho)$, $M_{\mathbf{k}}$ and $B_{\mathbf{k}}$, as defined by (4.43) and (4.44) on page 65, and by (4.70) on page 78, on the respective finite element spaces.

Typically, finite element matrices are assembled by applying specific quadrature rules. These quadrature rules are realized by evaluating the element shape functions at a finite set of quadrature points. During the matrix assembly, coefficients such as ρ are also evaluated at these quadrature points only. For convenience, we define the set

$$\mathcal{X}_h := \{\mathbf{x} \in \mathbb{R}^3 \mid \mathbf{x} \text{ is a quadrature point of } \mathcal{T}_h\}. \quad (7.4)$$

Since the coefficient ρ is only evaluated at the points in \mathcal{X}_h during the assembly of $\mathbf{A}_{h,\mathbf{k}}(\rho)$, we devise the following discretization schemes for the sets \mathcal{C} , $\bar{\mathcal{C}}$, and \mathcal{D} . First, we define the finite-dimensional, linear function space

$$\mathcal{E}_h := \{\xi_h \mid \xi_h : \mathcal{X}_h \rightarrow \mathbb{R}\}. \quad (7.5)$$

Clearly, the space \mathcal{E}_h consists of all real-valued functions, which are defined on the set of quadrature points of the finite element mesh \mathcal{T}_h . In the following, we shall use \mathcal{E}_h as a discretization of \mathcal{E} . An interpolation from the infinite-dimensional

Banach space \mathcal{E} into the discrete space \mathcal{E}_h can be defined as follows: Suppose that for every element T_i in the finite element mesh \mathcal{T}_h we are given an associated linear function space P_i that consists of real-valued polynomial functions defined on T_i . Then, we define the *interpolation operator* $\Pi_h : \mathcal{E} \rightarrow \mathcal{E}_h$ by

$$\Pi_h \xi := \pi(\xi)|_{\mathcal{X}_h} \quad \text{for all } \xi \in \mathcal{E}, \quad (7.6)$$

where $\pi(\xi)$ is the element-wise polynomial function that satisfies $\pi(\xi)|_{T_i} \in P_i$, as well as

$$\|\pi(\xi)|_{T_i} - \xi\|_{T_i, \infty} \leq \|p - \xi\|_{T_i, \infty} \quad \text{for all } p \in P_i, \quad i = 1, \dots, N.$$

Note that Π_h is a *Clément-type* interpolation operator (cf. [24]). In analogy to the function sets \mathcal{D} , \mathcal{C} , and $\bar{\mathcal{C}}$, we also define the discrete function sets

$$\mathcal{D}_h := \{\rho_h \in \mathcal{E}_h \mid \rho_h(\mathbf{x}) > 0 \text{ for all } \mathbf{x} \in \mathcal{X}_h\}, \quad (7.7)$$

$$\mathcal{C}_h := \{\rho_h \in \mathcal{E}_h \mid \rho_h(\mathbf{x}) \in \{\rho_{\min}, \rho_{\max}\} \text{ for all } \mathbf{x} \in \mathcal{X}_h\}, \quad (7.8)$$

$$\bar{\mathcal{C}}_h := \{\rho_h \in \mathcal{E}_h \mid \rho_{\min} \leq \rho_h(\mathbf{x}) \leq \rho_{\max} \text{ for all } \mathbf{x} \in \mathcal{X}_h\}. \quad (7.9)$$

We suppose that, for every index $j \in \mathbb{N}$, we have a function $\lambda_h : \mathcal{D}_h \times K_h \rightarrow \mathbb{R}$, such that for every coefficient $\rho \in \mathcal{D}$ and every vector $\mathbf{k} \in K_h$, the function value $\lambda_h(\Pi_h \rho, \mathbf{k})$ coincides with the j -th smallest eigenvalue of (7.3).

Given a function $\xi_h \in \mathcal{E}_h$, we define the integral of this function over Ω as

$$\int_{\Omega, h} \xi_h := \sum_{\mathbf{x} \in \mathcal{X}_h} w_{\mathbf{x}} \xi_h(\mathbf{x}), \quad (7.10)$$

where $w_{\mathbf{x}}$ denotes the quadrature weight corresponding to the quadrature point $\mathbf{x} \in \mathcal{X}_h$. Now, suppose that for every coefficient $\rho_h \in \mathcal{D}_h$ a subset $\Gamma_{J_h}(\rho_h)$ of \mathcal{E}_h can be computed, such that

$$\partial J_h(\rho_h) \subseteq \left\{ p \in \mathcal{E}_h^* \mid p(\eta_h) = \int_{\Omega, h} \gamma_h \eta_h, \quad \gamma_h \in \Gamma_{J_h}(\rho_h) \right\}. \quad (7.11)$$

Then, our aim is to construct a minimizing sequence $\rho_h^{(0)}, \rho_h^{(1)}, \rho_h^{(2)}, \dots$, etc. in \mathcal{C}_h for the discretized goal functional J_h . The initial coefficient $\rho^{(0)} \in \mathcal{C}_h$ is chosen arbitrarily. The coefficients $\rho_h^{(1)}, \rho_h^{(2)}, \dots$, are constructed iteratively according to

$$\rho_h^{(l)} := \rho_h^{(l-1)} - s_l \gamma_h^{(l-1)} \quad \text{for } l = 1, 2, \dots, \quad (7.12)$$

where $s_l > 0$ is a suitable step size, and where $\gamma_h^{(l-1)}$ is an element of the set $\Gamma_{J_h}(\rho_h^{(l-1)})$ for every $l \in \mathbb{N}$. Clearly, this approach corresponds to a generalized gradient method for minimization problem of the form (7.1) for $X = \mathcal{E}_h$, $f = J_h$, and for the inner product $\langle \cdot, \cdot \rangle_{\mathcal{E}_h}$ on \mathcal{E}_h , which is given by

$$\langle \xi_h, \eta_h \rangle_{\mathcal{E}_h} := \int_{\Omega, h} \xi_h \eta_h \quad \text{for all } \xi_h, \eta_h \in \mathcal{E}_h.$$

7.3 Incorporating Optimization Constraints

In the previous section we presented a general approach for an optimization algorithm than can be applied to discretized photonic band structure optimization problems (PBSOPs). In essence, this approach consists in constructing a minimizing sequence for the goal functional J_h by means of a generalized gradient method. The minimizing sequence is constructed iteratively according to (7.12). In each iteration a function in \mathcal{E}_h is added to the respective discretized coefficient.

So far, our approach does not account for certain optimization constraints which are inherent to PBSOPs. Recall that the discretized minimization problem (7.2) is posed on the admissible set \mathcal{C}_h , which consists of two-valued functions only. In order find an admissible, optimal density function, we have to impose an optimization constraint on the minimizing sequence $\rho_h^{(0)}, \rho_h^{(1)}, \rho_h^{(2)}, \dots$, etc., which ensures that every density function of that sequence is an admissible density function. Clearly, such an optimization constraint is given by requiring that

$$\rho_h^{(l)}(\mathbf{x}) \in \{\rho_{\min}, \rho_{\max}\} \quad \text{for all } \mathbf{x} \in \mathcal{X}_h, l = 0, 1, 2, \dots \quad (7.13)$$

Notice however, that by adding a function in \mathcal{E}_h to an admissible density function in \mathcal{C}_h , one generally obtains a function in $\mathcal{E}_h \setminus \mathcal{C}_h$. Therefore, we cannot expect that the optimization constraint (7.13) is automatically fulfilled by a standard generalized gradient method.

In the following we shall relax the optimization constraint on the minimizing sequence by considering the minimization problem

$$\underset{\rho_h \in \bar{\mathcal{C}}_h}{\text{minimize}} J_h(\rho_h).$$

The extended admissible set $\bar{\mathcal{C}}_h$ imposes the optimization constraint

$$\rho_{\min} \leq \rho_h^{(l)}(\mathbf{x}) \leq \rho_{\max} \quad \text{for all } \mathbf{x} \in \mathcal{X}_h, l = 0, 1, 2, \dots \quad (7.14)$$

on every minimizing sequence $\rho_h^{(0)}, \rho_h^{(1)}, \rho_h^{(2)}, \dots$, etc. that is constructed by a generalized gradient algorithm. Optimization constraints of the above form (7.14) are often referred to as *box constraints*. Again, one observes that adding a function in \mathcal{E}_h to a density function in $\bar{\mathcal{C}}_h$ generally yields a function in $\mathcal{E}_h \setminus \bar{\mathcal{C}}_h$. Therefore, one cannot expect that the general construction principle (7.12) yields a sequence in $\bar{\mathcal{C}}_h$.

In the following we present a method by which the box constraint (7.14) can be ensured. To this end, we define an unconstrained minimization problem on the finite-dimensional, linear function space \mathcal{E}_h . Suppose that $Y_h : \mathcal{E}_h \rightarrow \bar{\mathcal{C}}_h$ is an operator, which is given by

$$Y_h(\xi_h) := y \circ \xi_h \quad \text{for all } \xi_h \in \mathcal{E}_h, \quad (7.15)$$

where y is a strictly increasing, continuously differentiable function from \mathbb{R} onto the open interval $(\rho_{\min}, \rho_{\max})$. In the following, we shall refer to such an operator as a *funnel operator*, alluding to the fact that Y_h maps a function in $\xi_h \in \mathcal{E}_h$ to a function $Y_h(\xi_h) \in \bar{\mathcal{C}}_h$ by “funnelling” the function values of ξ_h from the “wider” interval $\mathbb{R} = (-\infty, \infty)$ into the “narrower” interval $(\rho_{\min}, \rho_{\max})$. Since the function y is assumed to be strictly increasing, we have that every funnel operator is a bijection.

Given a goal functional $J_h : \mathcal{D}_h \rightarrow \mathbb{R}$ and a specific funnel operator $Y_h : \mathcal{E}_h \rightarrow \bar{\mathcal{C}}_h$, we consider the minimization problem

$$\underset{\xi_h \in \mathcal{E}_h}{\text{minimize}} (J_h \circ Y_h)(\xi_h). \quad (7.16)$$

By constructing a minimizing sequence $\xi_h^{(0)}, \xi_h^{(1)}, \xi_h^{(2)}, \dots$, etc. in the function space \mathcal{E}_h for $J_h \circ Y_h$, we also obtain a minimizing sequence $\rho_h^{(0)}, \rho_h^{(1)}, \rho_h^{(2)}, \dots$, etc. in the extended admissible set $\bar{\mathcal{C}}_h$ for J_h . This minimizing sequence is given by

$$\rho_h^{(l)} := Y_h(\xi_h^{(l)}) \quad l = 0, 1, 2, \dots \quad (7.17)$$

Next, we consider the generalized gradients of the functional $J_h \circ Y_h$. One can easily show that the funnel operator Y_h is continuously Gâteaux differentiable in \mathcal{E}_h . At every point $\xi_h \in \mathcal{E}_h$, the corresponding Gâteaux differential $DY_h(\xi_h)$ of Y_h is given by

$$DY_h(\xi_h)[\eta_h] := (y' \circ \xi_h)\eta_h \quad \text{for all } \eta_h \in \mathcal{E}_h,$$

where y' denotes the derivative of the function y in (7.15). By Theorem 6.12 on page 117, we obtain the following inclusion relation for the generalized differentials of $J_h \circ Y_h$,

$$\partial(J_h \circ Y_h)(\xi_h) \subseteq \{p \circ DY_h(\xi_h) \mid p \in \partial J_h(Y_h(\xi_h))\} \quad \text{for all } \xi_h \in \mathcal{E}_h.$$

Furthermore, we obtain from (7.11) the inclusion relation

$$\partial(J_h \circ Y_h)(\xi_h) \subseteq \left\{ p \in \mathcal{E}_h^* \mid p(\eta_h) = \int_{\Omega, h} \gamma_h(y' \circ \xi_h)\eta_h, \gamma_h \in \Gamma_{J_h}(Y_h(\xi_h)) \right\}. \quad (7.18)$$

For convenience, we define for every function $\xi_h \in \mathcal{E}_h$ the set

$$\Gamma_{J_h \circ Y_h}(\xi_h) := \{(y' \circ \xi_h)\gamma_h \mid \gamma_h \in \Gamma_{J_h}(Y_h(\xi_h))\}.$$

With this, the inclusion relation (7.18) can be equivalently written as

$$\partial(J_h \circ Y_h)(\xi_h) \subseteq \left\{ p \in \mathcal{E}_h^* \mid p(\eta_h) = \int_{\Omega, h} \zeta_h \eta_h, \zeta_h \in \Gamma_{J_h \circ Y_h}(\xi_h) \right\}. \quad (7.19)$$

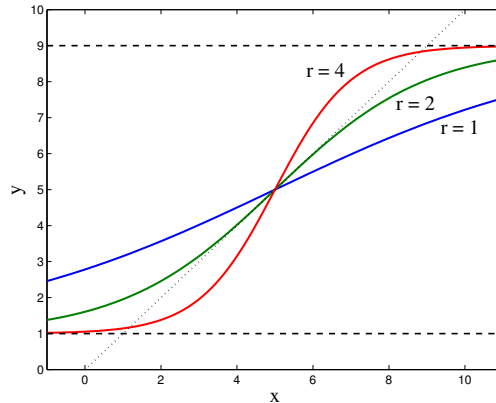


Figure 7.2: Graphs of the function y as defined by (7.20) for $\rho_{\min} = 1$, $\rho_{\max} = 9$ and various values for r . The dotted line indicates the diagonal $x = y$.

Notice that for every function $\xi_h \in \mathcal{E}_h$ the functions in the set $\Gamma_{J_h \circ Y_h}(\xi_h)$ are obtained from the functions in the set $\Gamma_{J_h}(Y_h(\xi_h))$ through a point-wise multiplication by the function $y' \circ \xi_h$. Therefore, computing the set $\Gamma_{J_h \circ Y_h}(\xi_h)$ is no more difficult than computing the set $\Gamma_{J_h}(Y_h(\xi_h))$. It should also be noted that $y' \circ \xi_h$ is a strictly positive function, since the function y is strictly increasing by assumption.

A minimizing sequence $\xi_h^{(0)}, \xi_h^{(1)}, \xi_h^{(2)}, \dots$, etc. in \mathcal{E}_h for the goal functional $J_h \circ Y_h$ can now be constructed using a generalized gradient method. After choosing an initial function $\xi_h^{(0)} \in \mathcal{E}_h$, the functions $\xi_h^{(1)}, \xi_h^{(2)}, \dots$, etc. are constructed according to

$$\xi_h^{(l)} := \xi_h^{(l-1)} + s_l \zeta_h^{(l-1)} \quad \text{for } l = 1, 2, \dots,$$

where $s_l > 0$ is a suitable step size, and where $\zeta_h^{(l-1)}$ is an element of $\Gamma_{J_h \circ Y_h}(\xi_h^{(l-1)})$, such that $-\zeta_h^{(l-1)}$ is a descent direction of $J_h \circ Y_h$ at the point $\xi_h^{(l-1)}$ for $l = 1, 2, \dots$, etc. From this sequence, one obtains through (7.17) a minimizing sequence $\rho_h^{(0)}, \rho_h^{(1)}, \rho_h^{(2)}, \dots$, etc. in the extended admissible set $\bar{\mathcal{C}}$ for the goal functional J_h .

We remark that the function $y : \mathbb{R} \rightarrow (\rho_{\min}, \rho_{\max})$ in (7.15) can be defined as

$$\begin{cases} y(x) := \frac{\rho_{\max} - \rho_{\min}}{1 - \exp(-t_r(x))} + \rho_{\min} & \text{for all } x \in \mathbb{R}, \\ t_r(x) := \frac{2r}{\rho_{\max} - \rho_{\min}}(x - \rho_{\min}) - r & \text{for all } x \in \mathbb{R} \end{cases} \quad (7.20)$$

for some number $r > 0$. In Figure 7.2 we depicted the function graphs of y as defined by (7.20) for various values of the parameter r . Notice that, for every $r > 0$,

the graph of y is a sigmoid curve, whose horizontal asymptotes are given by ρ_{\min} and ρ_{\max} . All function graphs intersect at the point $((\rho_{\min} + \rho_{\max})/2, (\rho_{\min} + \rho_{\max})/2)$, which is also the point of inflection of every graph. The larger the value of r , the steeper the function graph is at that point, and the faster it tends to the horizontal asymptotes. Finally, we remark that the function y defined in (7.20) is a so-called logistic function. It is also possible, however, to construct similar functions based on the inverse tangent function or the error function.

7.4 Preserving Symmetries

In each iteration of a generalized gradient method for the unconstrained minimization problem (7.16) the goal functional $J_h \circ Y_h$ is evaluated by computing an approximate photonic band structure. Typically, the band structure is computed only at a finite set of quasimomentum vectors $\mathbf{k}_1, \dots, \mathbf{k}_M$ that belong to the first Brillouin zone \mathbb{B} . It was discussed in Section 4.6, that under certain assumptions on the photonic crystal's point group the band structure is completely determined by its restriction to a so-called irreducible zone $K \subset \mathbb{B}$. The set of quasimomentum vectors $\mathbf{k}_1, \dots, \mathbf{k}_M$ is hence usually chosen from within this irreducible zone K . Typically, these quasimomentum vectors lie on the critical path in \mathbb{B} (see Figure 4.1(a) in Section 4.6).

It follows that a generalized gradient method has to guarantee, that the density functions of the minimizing sequence $\rho_h^{(0)}, \rho_h^{(1)}, \rho_h^{(2)}, \dots$, etc. defined by (7.17) all possess the same symmetry. This imposes a so-called *symmetry constraint* on the sequence. Provided that the point group of the initial photonic crystal's medium structure is isomorphic to a group G consisting of orthogonal transformations, the symmetry constraint is given by

$$\rho_h^{(l)} \circ \boldsymbol{\theta} = \rho_h^{(l)} \quad \text{for all } \boldsymbol{\theta} \in G, l = 0, 1, 2, \dots \quad (7.21)$$

In Section 2.4, we introduced the symmetrization operation $f \mapsto f^{(G)}$, where G is a finite group of isometries. This operation acts as an identity on the set of G -symmetric functions and maps each function f to a G -symmetric function $f^{(G)}$. We constitute that the symmetrization operation is realized on the finite-dimensional, linear function space \mathcal{E}_h by a so-called *symmetrization operator* $S_G : \mathcal{E}_h \rightarrow \mathcal{E}_h$ that is defined as

$$S_G(\xi_h) := \xi_h^{(G)} \quad \text{for all } \xi_h \in \mathcal{E}_h. \quad (7.22)$$

Here as in the following, we assume that the set of quadrature points \mathcal{X}_h is G -symmetric. From (2.11) on page 18 one easily deduces that S_G is a linear operator for every group G of orthogonal transformations. The action of S_G on the set \mathcal{E}_h can be realized as follows: Suppose that the set of quadrature points is given

by $\mathcal{X}_h = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$, where $N := |\mathcal{X}_h|$. Then, each function $\xi_h \in \mathcal{E}_h$ can be represented by a corresponding vector $\underline{\xi}_h \in \mathbb{R}^N$ by virtue of

$$\underline{\xi}_{h_i} = \xi_h(\mathbf{x}^{(i)}) \quad \text{for all } i = 1, \dots, N,$$

where $\underline{\xi}_{h_i}$ denotes the i -th component of the vector $\underline{\xi}_h$. Given a group G of orthogonal transformations, such that \mathcal{X}_h is G -symmetric, one can show that there exists a symmetric matrix $\mathbf{S} \in \mathbb{R}^{n \times n}$, such that

$$\underline{S}_G(\underline{\xi}_h) = \mathbf{S}\underline{\xi}_h \quad \text{for all } \xi_h \in \mathcal{E}_h.$$

As with the optimization constraint imposed by the extended admissible set $\overline{\mathcal{C}}_h$, we incorporate the symmetry constraint (7.21) into a generalized gradient method by considering yet another minimization problem, which is given by

$$\underset{\xi_h \in \mathcal{E}_h}{\text{minimize}} (J_h \circ Y_h \circ S_G)(\xi_h)$$

for some group G of orthogonal transformations. Since S_G is a linear operator, we have by Theorem 6.12 on page 117 that

$$\partial(J_h \circ Y_h \circ S_G)(\xi_h) \subseteq \left\{ p \circ S_G \mid p \in \partial(J_h \circ Y_h)(S_G(\xi_h)) \right\} \quad \text{for all } \xi_h \in \mathcal{E}_h.$$

From (7.18), we obtain the inclusion relation

$$\begin{aligned} \partial(J_h \circ Y_h \circ S_G)(\xi_h) \subseteq \left\{ p \in \mathcal{E}_h^* \mid p(\eta_h) = \int_{\Omega, h} \gamma_h(y' \circ \xi_h) S_G(\eta_h), \right. \\ \left. \gamma_h \in \Gamma_{J_h}((Y_h \circ S_G)(\xi_h)) \right\}. \end{aligned} \quad (7.23)$$

Recall that the integral over Ω in (7.18) is understood as the application of a quadrature rule on the set of quadrature point \mathcal{X}_h . Suppose that $\mathbf{W} \in \mathbb{R}^{N \times N}$ denotes the diagonal matrix with the corresponding quadrature weights on the diagonal. Given some arbitrary functions $\xi_h \in \mathcal{E}_h$ and $\gamma_h \in \Gamma_{J_h}(\xi_h)$, we find that

$$\begin{aligned} \int_{\Omega, h} \gamma_h(y' \circ \xi_h) S_G(\eta_h) &= (\underline{\gamma}_h \odot (\underline{y}' \circ \underline{\xi}_h)) \cdot \mathbf{W} \mathbf{S} \underline{\eta}_h \\ &= \left(\mathbf{S} (\underline{\gamma}_h \odot (\underline{y}' \circ \underline{\xi}_h)) \right) \cdot \mathbf{W} \underline{\eta}_h \\ &= \int_{\Omega, h} S_G(\gamma_h(y' \circ \xi_h)) \eta_h \quad \text{for all } \eta_h \in \mathcal{E}_h. \end{aligned}$$

Here, \odot denotes the *component-wise vector product*, which also known as the *Hadamard product* on \mathbb{R}^N . The above identity affords an alternative formulation

of the inclusion relation (7.23). In order to establish this formulation, we define for every function $\xi_h \in \mathcal{E}_h$ the set

$$\Gamma_{J_h \circ Y_h \circ S_G}(\xi_h) := \{S_G((y' \circ \xi_h)\gamma_h) \mid \gamma_h \in \Gamma_{J_h}((Y_h \circ S_G)(\xi_h))\}.$$

With this, the inclusion relation (7.23) can be equivalently written as

$$\partial(J_h \circ Y_h \circ S_G)(\xi_h) \subseteq \{p \in \mathcal{E}_h^* \mid p(\eta_h) = \int_{\Omega, h} v_h \eta_h, v_h \in \Gamma_{J_h \circ Y_h \circ S_G}(\xi_h)\}.$$

From this we deduce that the set $\Gamma_{J_h \circ Y_h \circ S_G}(\xi_h)$ is a superset of the generalized gradient of $J_h \circ Y_h \circ S_G$ at ξ_h for every function $\xi_h \in \mathcal{E}_h$.

We can now proceed as in the previous section and define a generalized gradient method, which constructs a minimizing sequence $\xi_h^{(0)}, \xi_h^{(1)}, \xi_h^{(2)}, \dots$, etc. in \mathcal{E}_h for the goal functional $J_h \circ Y_h \circ S_G$. In the l -th iteration of this method, the descent direction $-v_h^{(l-1)}$ is chosen such that $v_h^{(l-1)} \in \Gamma_{J_h \circ Y_h \circ S_G}(\xi_h^{(l-1)})$ for $l = 1, 2, \dots$, etc. By construction, every function in the set $\Gamma_{J_h \circ Y_h \circ S_G}(\xi_h^{(l-1)})$ is G -symmetric. Provided that the initial function $\xi_h^{(0)} \in \mathcal{E}_h$ is also G -symmetric, we thus obtain a minimizing sequence, which consists of G -symmetric functions in \mathcal{E}_h . In particular, we have that

$$S_G(\xi_h^{(l)}) = \xi_h^{(l)} \quad \text{for } l = 0, 1, 2, \dots$$

From the sequence $\xi_h^{(0)}, \xi_h^{(1)}, \xi_h^{(2)}, \dots$, we obtain by (7.17) a minimizing sequence $\rho_h^{(0)}, \rho_h^{(1)}, \rho_h^{(2)}, \dots$, etc. in $\overline{\mathcal{C}}_h$ for the goal functional J_h , such that every density function of that sequence is G -symmetric.

7.5 The Algorithm

Based on the approaches that were introduced in the previous Section 7.3 and Section 7.4, we implemented an optimization algorithm for band structure optimization problems. In this section we describe this algorithm and comment on some implementation details. The algorithm is listed in pseudo-code as Algorithm 7.1. It was implemented in MATLAB as well as in C++ using the parallel finite element library M++ (see [76], [77]). In the following we shall refer to the functions $\xi_h^{(l)}$ in Algorithm 7.1 as *free density functions*.

The optimization algorithm constitutes a generalized gradient method for the locally Lipschitz minimization problem

$$\underset{\xi_h \in \mathcal{E}_h}{\text{minimize}} (J_h \circ Y_h \circ S_G)(\xi_h),$$

as introduced in Section 7.1. A minimizing sequence $\xi_h^{(0)}, \xi_h^{(1)}, \xi_h^{(2)}, \dots$, etc. of free density functions in \mathcal{E}_h is iteratively constructed from a given initial free

Algorithm 7.1 The Generalized Gradient Algorithm**Input:**

$J_h : \mathcal{D}_h \rightarrow \mathbb{R}$	<i>Goal functional</i>
G	<i>Group of orthogonal transformations (Point group)</i>
$\xi_h^{(0)}$	<i>Initial free density</i>

Parameters:

$s_{\min} > 0$	<i>Minimal step size</i>
$s_{\max} > s_{\min}$	<i>Maximal step size</i>
$\Delta_{\min} > 0$	<i>Minimal decrease</i>

Algorithm:

```

 $\rho_h^{(0)} := Y_h(\xi_h^{(0)})$ 
compute  $z_0 := J_h(\rho_h^{(0)})$ 
 $s_0 := s_{\max}$ 
for  $l = 1, 2, \dots$  do
  compute  $\Gamma_{J_h \circ Y_h \circ S_G}(\xi_h^{(l-1)})$ 
  choose a descent direction  $-\tilde{v}_h^{(l-1)} \in -\Gamma_{J_h \circ Y_h \circ S_G}(\xi_h^{(l-1)})$ 
   $v_h^{(l-1)} := \tilde{v}_h^{(l-1)} / \|\tilde{v}_h^{(l-1)}\|_\infty$ 
   $s_{l,1} := \min\{2s_{l-1}, s_{\max}\}$ 
  for  $p = 1, 2, \dots$  do
     $\xi_h^{(l,p)} := \xi_h^{(l-1)} - s_{l,p} v_h^{(l-1)}$ 
     $\rho_h^{(l,p)} := Y_h(\xi_h^{(l,p)})$ 
    compute  $z_{l,p} := J_h(\rho_h^{(l,p)})$ 
    if  $z_{l-1} - z_{l,p} > \Delta_{\min}$  then
       $\xi_h^{(l)} := \xi_h^{(l,p)}$ 
       $z_l := z_{l,p}$ 
       $t_l := s_{l,p}$ 
      break
    end if
     $s_{l,p+1} := s_{l,p}/2$ 
    if  $s_{l,p+1} < s_{\min}$  then
      return  $\rho_h^{(l-1)}$ 
    end if
  end for
end for

```

function $\xi_h^{(0)} \in \mathcal{E}_h$ according to

$$\xi_h^{(l)} := \xi_h^{(l-1)} - s_l v_h^{(l-1)} \quad \text{for } l = 1, 2, \dots \quad (7.24)$$

In every iteration, the function $v_h^{(l-1)}$ is chosen from the set $\Gamma_{J_h \circ Y_h \circ S_G}(\xi_h^{(l-1)})$, such

that $-v_h^{(l-1)}$ is a descent direction of $J_h \circ Y_h \circ S_G$ at $\xi_h^{(l-1)}$, where $l = 1, 2, \dots$. The choice of the descent direction is accomplished by an appropriate choice strategy. In general, one has to devise a separate choice strategy for every goal functional J_h . In the following section we briefly discuss this issue and present a choice strategy, which is appropriate for the gap width functionals defined in Section 5.5.

The step sizes $s_l > 0$ are chosen adaptively using a trial-and-error scheme. Starting with a given maximal step size s_{\max} the step size is divided by 2, each time an update of the current free density function according to (7.24) fails to yield a prescribed minimal decrease Δ_{\min} in the goal value. In order to prevent excessively small step sizes, the algorithm always doubles the step size before attempting an update by a new descent direction. The algorithm terminates, when the step size becomes smaller than a prescribed minimal step size $s_{\min} > 0$. In this case, we shall say that the algorithm has *converged*. The *final free density function* is given by the last free density function $\xi_h^{(L)} \in \mathcal{E}_h$, for which a decrease of the goal value could be obtained. The final free density defines a corresponding *final density function* $\rho_h^{(L)} \in \bar{\mathcal{C}}_h$, which is given by

$$\rho_h^{(L)} := Y_h(\xi_h^{(L)}).$$

This final density function $\rho_h^{(L)}$ is returned by the algorithm.

By construction the final density function $\rho_h^{(L)}$ is G -symmetric. Due to the definition of the funnel operator Y_h , however, $\rho_h^{(L)}$ is always a function in $\bar{\mathcal{C}}_h \setminus \mathcal{C}_h$. Hence, the optimization algorithm is not capable to return admissible density functions. In numerical experiments we found, however, that the final density function was almost two-valued, attaining function values close to ρ_{\min} and ρ_{\max} at most quadrature points. We remark that an admissible function can always be obtained from a final density function $\rho_h^{(L)}$, by applying the so-called *threshold operator* $H_h : \bar{\mathcal{C}}_h \rightarrow \mathcal{C}_h$ to it. The threshold operator is defined by

$$H_h(\rho_h)(\mathbf{x}) := \begin{cases} \rho_{\min} & \text{if } \rho_h(\mathbf{x}) < \frac{\rho_{\min} + \rho_{\max}}{2}, \\ \rho_{\max} & \text{else} \end{cases} \quad \text{for all } \mathbf{x} \in \mathcal{X}_h. \quad (7.25)$$

and for all $\rho_h \in \bar{\mathcal{C}}_h$. We used Algorithm 7.1 to solve several photonic band gap maximization problems. Some of the results are presented in Chapter 9.

7.6 Choosing Descent Directions

A crucial step each iteration of the generalized gradient method, which is given by Algorithm 7.1, is the choice of the descent directions $-v_h^{(l-1)}$. As we mentioned in the previous section, this choice is usually accomplished through a specific choice

strategy. In this section we briefly comment on possible choice strategies and develop a specific choice strategy for the band gap functionals defined in Section 5.5.

In certain situations, the choice of a descent direction is trivial. Suppose, for example, that the set $\Gamma_{J_h}(\xi_h)$ is a singleton. Then, the set $\Gamma_{J_h \circ Y_h \circ S_G}(\xi_h^{(l-1)})$ is also a singleton and the functional $J_h \circ Y_h \circ S_G$ is hence strictly differentiable at $\xi_h^{(l-1)}$. In this case the set $-\Gamma_{J_h \circ Y_h \circ S_G}(\xi_h^{(l-1)})$ contains the negative strict gradient of $J_h \circ Y_h \circ S_G$ at $\xi_h^{(l-1)}$ as the only element, which constitutes a descent direction. Furthermore, when $-\Gamma_{J_h \circ Y_h \circ S_G}(\xi_h^{(l-1)})$ is known to be equal to the negative of the generalized gradient of $J_h \circ Y_h \circ S_G$ at $\xi_h^{(l-1)}$, one can choose the element with minimal norm as a descent direction according to Lemma 7.1. In the general case, however, the set $\Gamma_{J_h \circ Y_h \circ S_G}(\xi_h^{(l-1)})$ is strictly larger than the generalized gradient.

We now turn to the specific case, where the goal functional J_h is given by a *discretized gap width functional* $J_{j,h} : \mathcal{D}_h \rightarrow \mathbb{R}$ for some fixed index $j \in \mathbb{N}$. We assume that the $J_{j,h}$ is given by

$$J_{j,h}(\rho_h) := \max_{\mathbf{k} \in K_h} \omega_{j,h}(\rho_h, \mathbf{k}) - \min_{\mathbf{k} \in K_h} \omega_{j+1,h}(\rho_h, \mathbf{k}) \quad \text{for all } \rho_h \in \mathcal{D}_h,$$

where K_h denotes a finite set of quasimomentum vectors in the first Brillouin zone \mathbb{B} , and where $\omega_{j,h}(\rho_h, \mathbf{k})$ denotes the square-root of the j -th smallest eigenvalue $\lambda_{j,h}(\rho_h, \mathbf{k})$ of the generalized matrix eigenvalue problem (7.3) on page 131 for all $\rho_h \in \mathcal{D}_h$ and all $\mathbf{k} \in K_h$. We also define the *discretized band edge functionals* $J_{\text{lo},j,h}, J_{\text{up},j,h} : \mathcal{D}_h \rightarrow \mathbb{R}$ by

$$\begin{aligned} J_{\text{lo},j,h}(\rho_h) &:= \max_{\mathbf{k} \in K_h} \omega_{j,h}(\rho_h, \mathbf{k}) && \text{for all } \rho_h \in \mathcal{D}_h, \\ J_{\text{up},j,h}(\rho_h) &:= \min_{\mathbf{k} \in K_h} \omega_{j+1,h}(\rho_h, \mathbf{k}), && \text{for all } \rho_h \in \mathcal{D}_h. \end{aligned}$$

In the following we assume that every eigenvalue $\lambda_{j,h}(\rho_h, \mathbf{k})$ is simple. We remark that this assumption is usually fulfilled due to numerical errors. Then, by a similar result to that of Theorem 6.16 on page 121, we then have that

$$\partial \lambda_{j,h}(\cdot, \mathbf{k})(\rho_h) = \left\{ p \in \mathcal{E}_h^* \mid p(\eta_h) := \int_{\Omega,h} \eta |\mathcal{I}_h[(\nabla + i\mathbf{k})_h \times \mathbf{u}_{j,h}(\rho_h, \mathbf{k})]|^2 \right\},$$

where $\mathbf{u}_{j,h}(\rho_h, \mathbf{k})$ denotes an eigenvalue corresponding to $\lambda_{j,h}(\rho_h, \mathbf{k})$, and where $(\nabla + i\mathbf{k})_h \times$ denotes a modified curl operator, which is defined appropriately on the finite element space \mathbf{W}_h . Recall that \mathcal{I}_h denotes the interpolation operator with image space \mathcal{E}_h (see Section 7.2).

Since the generalized differential $\partial \lambda_{j,h}(\cdot, \mathbf{k})(\rho_h)$ is a singleton, we have by Proposition 6.7(d) on page 114 that the functional $\lambda_{j,h}(\cdot, \mathbf{k})$ is strictly differentiable in a neighbourhood of ρ_h for every vector $\mathbf{k} \in K_h$.

We further assume that

$$\arg \max_{K_h} \lambda_{j,h}(\rho_{j,h}, \cdot) = \{\mathbf{k}_{\text{lo}}^{(1)}, \dots, \mathbf{k}_{\text{lo}}^{(N_{\text{lo}})}\}, \quad (7.26)$$

$$\arg \min_{K_h} \lambda_{j+1,h}(\rho_{j,h}, \cdot) = \{\mathbf{k}_{\text{up}}^{(1)}, \dots, \mathbf{k}_{\text{up}}^{(N_{\text{up}})}\} \quad (7.27)$$

for small numbers $N_{\text{lo}}, N_{\text{up}} \in \mathbb{N}$. For convenience, we define the functionals $\tilde{\omega}_{\text{lo}}^{(i)}, \tilde{\omega}_{\text{up}}^{(i)} : \mathcal{E}_h \rightarrow \mathbb{R}$ by

$$\begin{aligned} \tilde{\omega}_{\text{lo}}^{(i)} &:= \omega_{j,h}(\cdot, \mathbf{k}_{\text{lo}}^{(i)}) \circ Y_h \circ S_G & i = 1, \dots, N_{\text{lo}}, \\ \tilde{\omega}_{\text{up}}^{(i)} &:= \omega_{j+1,h}(\cdot, \mathbf{k}_{\text{up}}^{(i)}) \circ Y_h \circ S_G & i = 1, \dots, N_{\text{up}}, \end{aligned}$$

where Y_h and S_G are the funnel operator (see (7.15) in Section 7.3) and the symmetrization operator (see (7.22) in Section 7.4) in Algorithm 7.1. Given a fixed function $\xi_h \in \mathcal{E}_h$, we define the functions $v_{\text{lo}}^{(i)}, v_{\text{up}}^{(i)} \in \mathcal{E}_h$ by

$$\begin{aligned} v_{\text{lo}}^{(i)} &:= \frac{S_G\left((y' \circ \xi_h) \left| \mathcal{I}_h[(\nabla + i\mathbf{k}_{\text{lo}}^{(i)})_h \times \mathbf{u}_{j,h}(\rho_h, \mathbf{k}_{\text{lo}}^{(i)})] \right|^2\right)}{2\omega_{j,h}(\rho_h, \kappa_{\text{lo}}^{(i)})} & i = 1, \dots, N_{\text{lo}}, \\ v_{\text{up}}^{(i)} &:= \frac{S_G\left((y' \circ \xi_h) \left| \mathcal{I}_h[(\nabla + i\mathbf{k}_{\text{up}}^{(i)})_h \times \mathbf{u}_{j+1,h}(\rho_h, \mathbf{k})] \right|^2\right)}{2\omega_{j+1,h}(\rho_h, \kappa_{\text{up}}^{(i)})} & i = 1, \dots, N_{\text{up}}. \end{aligned}$$

Then, we have that generalized differentials of the functionals $\tilde{\omega}_{\text{lo}}^{(i)}$ and $\tilde{\omega}_{\text{up}}^{(i)}$ at ξ_h are given by

$$\begin{aligned} \partial \tilde{\omega}_{\text{lo}}^{(i)}(\xi_h) &= \left\{ p \in \mathcal{E}_h^* \left| p(\eta_h) := \int_{\Omega,h} v_{\text{lo}}^{(i)} \eta_h \right. \right\} & i = 1, \dots, N_{\text{lo}}, \\ \partial \tilde{\omega}_{\text{up}}^{(i)}(\xi_h) &= \left\{ p \in \mathcal{E}_h^* \left| p(\eta_h) := \int_{\Omega,h} v_{\text{up}}^{(i)} \eta_h \right. \right\} & i = 1, \dots, N_{\text{up}}. \end{aligned}$$

Since the generalized differentials $\partial \tilde{\omega}_{\text{lo}}^{(i)}(\xi_h)$ are singletons, we have that the functional functionals $\tilde{\omega}_{\text{lo}}^{(i)}$ are strictly differentiable at ξ_h for all $i = 1, \dots, N_{\text{lo}}$. By the same argument, we have that the functionals $\tilde{\omega}_{\text{up}}^{(i)}$ are also strictly differentiable at ξ_h for all $i = 1, \dots, N_{\text{up}}$. In particular, we have that

$$\tilde{\omega}_{\text{lo}}^{(i)}(\xi_h + \eta_h) = \tilde{\omega}_{\text{lo}}^{(i)}(\xi_h) + \int_{\Omega,h} v_{\text{lo}}^{(i)} \eta_h + o(\|\eta\|_{\mathcal{E}_h}) \quad \text{as } \eta_h \rightarrow 0 \quad i = 1, \dots, N_{\text{lo}}, \quad (7.28)$$

$$\tilde{\omega}_{\text{up}}^{(i)}(\xi_h + \eta_h) = \tilde{\omega}_{\text{up}}^{(i)}(\xi_h) + \int_{\Omega,h} v_{\text{up}}^{(i)} \eta_h + o(\|\eta\|_{\mathcal{E}_h}) \quad \text{as } \eta_h \rightarrow 0 \quad i = 1, \dots, N_{\text{up}}. \quad (7.29)$$

For convenience we define the functionals $\tilde{J}, \tilde{J}_{\text{lo}}, \tilde{J}_{\text{up}} : \mathcal{E}_h \rightarrow \mathbb{R}$ by

$$\begin{aligned}\tilde{J} &:= J_{j,h} \circ Y_h \circ S_G, \\ \tilde{J}_{\text{lo}} &:= J_{\text{lo},j,h} \circ Y_h \circ S_G, \\ \tilde{J}_{\text{up}} &:= J_{\text{up},j,h} \circ Y_h \circ S_G.\end{aligned}$$

Clearly, we then have that

$$\tilde{J} = \tilde{J}_{\text{lo}} - \tilde{J}_{\text{up}}.$$

We also define the sets $\tilde{\Gamma}_{\text{lo}}, \tilde{\Gamma}_{\text{up}}, \tilde{\Gamma} \in \mathcal{E}_h$ as

$$\begin{aligned}\tilde{\Gamma}_{\text{lo}} &:= \overline{\text{conv}}\{v_{\text{up}}^{(1)}, \dots, v_{\text{up}}^{(N_{\text{lo}})}\}, \\ \tilde{\Gamma}_{\text{up}} &:= \overline{\text{conv}}\{v_{\text{up}}^{(1)}, \dots, v_{\text{up}}^{(N_{\text{up}})}\}, \\ \tilde{\Gamma} &:= \tilde{\Gamma}_{\text{lo}} - \tilde{\Gamma}_{\text{up}}.\end{aligned}$$

According to (7.26)–(7.27) and Proposition 6.15 on page 119, we have that

$$\begin{aligned}\partial\tilde{J}_{\text{lo}}(\xi_h) &\subseteq \left\{ p \in \mathcal{E}_h^* \mid p(\eta_h) = \int_{\Omega,h} v_h \eta_h, v_h \in \tilde{\Gamma}_{\text{lo}} \right\} \\ \partial\tilde{J}_{\text{up}}(\xi_h) &\subseteq \left\{ p \in \mathcal{E}_h^* \mid p(\eta_h) = \int_{\Omega,h} v_h \eta_h, v_h \in \tilde{\Gamma}_{\text{up}} \right\}.\end{aligned}$$

Hence, we obtain by Proposition 6.9 and Proposition 6.10 on page 116 the following inclusion relation for the generalized differential of \tilde{J} at ξ_h ,

$$\partial\tilde{J}(\xi_h) \subseteq \left\{ p \in \mathcal{E}_h^* \mid p(\eta_h) = \int_{\Omega,h} v_h \eta_h, v_h \in \tilde{\Gamma} \right\}.$$

From this, we deduce that the set $\tilde{\Gamma}$ is a superset of the generalized gradient of \tilde{J} at ξ_h . We remark that the set $\tilde{\Gamma}$ reduces to a singleton, if and only if $N_{\text{lo}} = N_{\text{up}} = 1$.

Our aim now is to indentify a descent direction $-v_h \in -\tilde{\Gamma}$ for the goal functional \tilde{J} at ξ_h . To this end we derive a linear program, whose solution determines a good candidate for such a descent direction. Note that by definition of the functionals $\tilde{\omega}_{\text{lo}}^{(i)}$ and $\tilde{\omega}_{\text{up}}^{(i)}$ we have that

$$\begin{aligned}\tilde{J}_{\text{lo}}(\xi_h) &= \tilde{\omega}_{\text{lo}}^{(1)}(\xi_h) = \dots = \tilde{\omega}_{\text{lo}}^{(N_{\text{lo}})}(\xi_h), \\ \tilde{J}_{\text{up}}(\xi_h) &= \tilde{\omega}_{\text{up}}^{(1)}(\xi_h) = \dots = \tilde{\omega}_{\text{up}}^{(N_{\text{up}})}(\xi_h).\end{aligned}$$

Under the assumption that the function identities

$$\tilde{J}_{\text{lo}} = \max\{\tilde{\omega}_{\text{lo}}^{(1)}, \dots, \tilde{\omega}_{\text{lo}}^{(N_{\text{lo}})}\},$$

$$\tilde{J}_{\text{up}} = \min\{\tilde{\omega}_{\text{up}}^{(1)}, \dots, \tilde{\omega}_{\text{up}}^{(N_{\text{up}})}\}$$

hold in some neighbourhood of ξ_h , we obtain by (7.28) and (7.29) the following asymptotic behaviour of \tilde{J} at ξ_h ,

$$\begin{aligned} & \tilde{J}(\xi_h + \eta_h) - \tilde{J}(\xi_h) \\ &= \tilde{J}_{\text{lo}}(\xi_h + \eta_h) - \tilde{J}_{\text{lo}}(P(\xi)) - \tilde{J}_{\text{up}}(\xi_h + \eta_h) + \tilde{J}_{\text{up}}(\xi_h) \\ &= \max\{\tilde{\omega}_{\text{lo}}^{(1)}(\xi_h + \eta_h) - \tilde{\omega}_{\text{lo}}^{(1)}(\xi_h), \dots, \tilde{\omega}_{\text{lo}}^{(N_{\text{lo}})}(\xi_h + \eta_h) - \tilde{\omega}_{\text{lo}}^{(N_{\text{lo}})}(\xi_h)\} \\ &\quad - \min\{\tilde{\omega}_{\text{up}}^{(1)}(\xi_h + \eta_h) - \tilde{\omega}_{\text{up}}^{(1)}(\xi_h), \dots, \tilde{\omega}_{\text{up}}^{(N_{\text{up}})}(\xi_h + \eta_h) - \tilde{\omega}_{\text{up}}^{(N_{\text{up}})}(\xi_h)\} \\ &= \max\left\{\int_{\Omega, h} v_{\text{lo}}^{(1)} \eta_h, \dots, \int_{\Omega, h} v_{\text{lo}}^{(N_{\text{lo}})} \eta_h\right\} - \min\left\{\int_{\Omega, h} v_{\text{up}}^{(1)} \eta_h, \dots, \int_{\Omega, h} v_{\text{up}}^{(N_{\text{up}})} \eta_h\right\} \\ &\quad + o(\|\eta_h\|_{\mathcal{E}_h}) \quad \text{as } \eta_h \rightarrow 0. \end{aligned}$$

This asymptotic behaviour implies, that a direction $\eta_h \in \mathcal{E}_h$, which minimizes

$$\max\left\{\int_{\Omega, h} v_{\text{lo}}^{(1)} \eta_h, \dots, \int_{\Omega, h} v_{\text{lo}}^{(N_{\text{lo}})} \eta_h\right\} - \min\left\{\int_{\Omega, h} v_{\text{up}}^{(1)} \eta_h, \dots, \int_{\Omega, h} v_{\text{up}}^{(N_{\text{up}})} \eta_h\right\}$$

is a good candidate for a descent direction of \tilde{J} at ξ_h . Now, suppose that the direction η_h is chosen from the set $-\tilde{\Gamma}$. Then, by definition of the sets $\tilde{\Gamma}$, $\tilde{\Gamma}_{\text{lo}}$ and $\tilde{\Gamma}_{\text{up}}$, there exists non-negative numbers $\theta_{\text{lo}}^{(1)}, \dots, \theta_{\text{lo}}^{(N_{\text{lo}})}, \theta_{\text{up}}^{(1)}, \dots, \theta_{\text{up}}^{(N_{\text{up}})} \in \mathbb{R}$, such that

$$0 \leq \theta_{\text{lo}}^{(i)} \leq 1 \quad i = 1, \dots, N_{\text{lo}}, \quad (7.30)$$

$$0 \leq \theta_{\text{up}}^{(i)} \leq 1 \quad i = 1, \dots, N_{\text{up}}, \quad (7.31)$$

$$\sum_{i=1}^{N_{\text{lo}}} \theta_{\text{lo}}^{(i)} = 1, \quad (7.32)$$

$$\sum_{i=1}^{N_{\text{up}}} \theta_{\text{up}}^{(i)} = 1 \quad (7.33)$$

and such that

$$\eta_h = - \sum_{\ell=1}^{N_{\text{lo}}} \theta_{\text{lo}}^{(\ell)} v_{\text{lo}}^{(\ell)} + \sum_{\ell=1}^{N_{\text{up}}} \theta_{\text{up}}^{(\ell)} v_{\text{up}}^{(\ell)}. \quad (7.34)$$

It follows that

$$\int_{\Omega,h} v_{\text{lo}}^{(i)} \eta_h = - \sum_{i=1}^{N_{\text{lo}}} \theta_{\text{lo}}^{(\ell)} \int_{\Omega,h} v_{\text{lo}}^{(i)} v_{\text{lo}}^{(\ell)} + \sum_{i=1}^{N_{\text{up}}} \theta_{\text{up}}^{(\ell)} \int_{\Omega,h} v_{\text{lo}}^{(i)} v_{\text{up}}^{(\ell)}, \quad i = 1, \dots, N_{\text{lo}}, \quad (7.35)$$

$$\int_{\Omega,h} v_{\text{up}}^{(i)} \eta_h = - \sum_{i=1}^{N_{\text{lo}}} \theta_{\text{lo}}^{(\ell)} \int_{\Omega,h} v_{\text{up}}^{(i)} v_{\text{lo}}^{(\ell)} + \sum_{i=1}^{N_{\text{up}}} \theta_{\text{up}}^{(\ell)} \int_{\Omega,h} v_{\text{up}}^{(i)} v_{\text{up}}^{(\ell)}, \quad i = 1, \dots, N_{\text{up}}. \quad (7.36)$$

We can now formulate a minimization problem, which is given by a linear program. To this end we define the functions $f_{\text{lo}}, f_{\text{up}} : -\tilde{\Gamma} \rightarrow \mathbb{R}$ by

$$f_{\text{lo}}(\eta_h) := \max \left\{ \int_{\Omega,h} v_{\text{lo}}^{(1)} \eta_h, \dots, \int_{\Omega,h} v_{\text{lo}}^{(N_{\text{lo}})} \eta_h \right\} \quad \text{for all } \eta_h \in -\tilde{\Gamma},$$

$$f_{\text{up}}(\eta_h) := \min \left\{ \int_{\Omega,h} v_{\text{up}}^{(1)} \eta_h, \dots, \int_{\Omega,h} v_{\text{up}}^{(N_{\text{up}})} \eta_h \right\} \quad \text{for all } \eta_h \in -\tilde{\Gamma}$$

and consider the minimization problem

$$\text{minimize}_{\eta_h \in -\tilde{\Gamma}} (f_{\text{lo}}(\eta_h) - f_{\text{up}}(\eta_h)). \quad (7.37)$$

By introducing the formal variables ζ_{lo} and ζ_{up} for $f_{\text{lo}}(\eta_h)$ and for $f_{\text{up}}(\eta_h)$, we obtain from (7.35)–(7.36) the constraints

$$- \sum_{\ell=1}^{N_{\text{lo}}} \theta_{\text{lo}}^{(\ell)} \int_{\Omega,h} v_{\text{lo}}^{(i)} v_{\text{lo}}^{(\ell)} + \sum_{\ell=1}^{N_{\text{up}}} \theta_{\text{up}}^{(\ell)} \int_{\Omega,h} v_{\text{lo}}^{(i)} v_{\text{up}}^{(\ell)} - \zeta_{\text{lo}} \leq 0 \quad \ell = 1, \dots, N_{\text{lo}},$$

$$\sum_{\ell=1}^{N_{\text{lo}}} \theta_{\text{lo}}^{(\ell)} \int_{\Omega,h} v_{\text{lo}}^{(i)} v_{\text{lo}}^{(\ell)} - \sum_{\ell=1}^{N_{\text{up}}} \theta_{\text{up}}^{(\ell)} \int_{\Omega,h} v_{\text{lo}}^{(i)} v_{\text{up}}^{(\ell)} + \zeta_{\text{up}} \leq 0 \quad \ell = 1, \dots, N_{\text{up}}.$$

Recall that the coefficients $\theta_{\text{lo}}^{(i)}$ and $\theta_{\text{up}}^{(i)}$ have to obey the constraints (7.30)–(7.33). By defining the unknown vector $\mathbf{x} \in \mathbb{R}^{N_{\text{lo}}+N_{\text{up}}+2}$ as

$$\mathbf{x} := \left(\theta_{\text{lo}}^{(1)}, \dots, \theta_{\text{lo}}^{(N_{\text{lo}})}, \theta_{\text{up}}^{(1)}, \dots, \theta_{\text{up}}^{(N_{\text{up}})}, \zeta_{\text{lo}}, \zeta_{\text{up}} \right)^{\text{T}}, \quad (7.38)$$

the minimization problem (7.37) can be formulated as a linear program according to

$$\begin{aligned} & \text{minimize}_{\mathbf{x}} \mathbf{f} \cdot \mathbf{x} \\ & \text{subject to } \mathbf{A}\mathbf{x} \leq \mathbf{0}, \quad \mathbf{C}\mathbf{x} = \mathbf{d}, \quad \mathbf{l} \leq \mathbf{x} \leq \mathbf{u}. \end{aligned}$$

The vector $\mathbf{f} \in \mathbb{R}^{N^{\text{lo}}+N^{\text{up}}+2}$ and the vectors $\mathbf{l}, \mathbf{u} \in \overline{\mathbb{R}}^{N^{\text{lo}}+N^{\text{up}}+2}$ are given by

$$\mathbf{f} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ -1 \end{pmatrix}, \quad \mathbf{l} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 0 \\ \vdots \\ 0 \\ -\infty \\ -\infty \end{pmatrix}, \quad \mathbf{u} = \begin{pmatrix} 1 \\ \vdots \\ 1 \\ 1 \\ \vdots \\ 1 \\ \infty \\ \infty \end{pmatrix}.$$

The matrix $\mathbf{A} \in \mathbb{R}^{(N^{\text{lo}}+N^{\text{up}}) \times (N^{\text{lo}}+N^{\text{up}}+2)}$ is given as a block matrix by

$$\mathbf{A} := \begin{pmatrix} -\mathbf{A}_{\text{lo,lo}} & \mathbf{A}_{\text{lo,up}} & -\mathbf{1}_{N^{\text{lo}},1} & \mathbf{0}_{N^{\text{lo}},1} \\ \mathbf{A}_{\text{up,lo}} & -\mathbf{A}_{\text{up,up}} & \mathbf{0}_{N^{\text{up}},1} & \mathbf{1}_{N^{\text{up}},1} \end{pmatrix}.$$

Here, we denote by $\mathbf{0}_{m,n}$ and $\mathbf{1}_{m,n}$ for given numbers $m, n \in \mathbb{N}$ those $m \times n$ matrices with components equal to zero and one, respectively. The matrices $\mathbf{A}_{\text{lo,lo}} \in \mathbb{R}_{N^{\text{up}} \times N^{\text{up}}}$, $\mathbf{A}_{\text{lo,up}} \in \mathbb{R}^{N^{\text{up}} \times N^{\text{up}}}$, $\mathbf{A}_{\text{up,lo}} \in \mathbb{R}^{N^{\text{up}} \times N^{\text{up}}}$ and $\mathbf{A}_{\text{up,up}} \in \mathbb{R}^{N^{\text{up}} \times N^{\text{up}}}$ are given by

$$\begin{aligned} [\mathbf{A}_{\text{lo,lo}}]_{i,\ell} &= \int_{\Omega,h} v_{\text{lo}}^{(i)} v_{\text{lo}}^{(\ell)} & i = 1, \dots, N^{\text{lo}}, \ell = 1, \dots, N^{\text{lo}}, \\ [\mathbf{A}_{\text{lo,up}}]_{i,\ell} &= \int_{\Omega,h} v_{\text{lo}}^{(i)} v_{\text{up}}^{(\ell)} & i = 1, \dots, N^{\text{lo}}, \ell = 1, \dots, N^{\text{up}}, \\ [\mathbf{A}_{\text{up,lo}}]_{i,\ell} &= \int_{\Omega,h} v_{\text{up}}^{(i)} v_{\text{lo}}^{(\ell)} & i = 1, \dots, N^{\text{up}}, \ell = 1, \dots, N^{\text{lo}}, \\ [\mathbf{A}_{\text{up,up}}]_{i,\ell} &= \int_{\Omega,h} v_{\text{up}}^{(i)} v_{\text{up}}^{(\ell)} & i = 1, \dots, N^{\text{up}}, \ell = 1, \dots, N^{\text{up}}. \end{aligned}$$

Finally, the matrix $\mathbf{C} \in \mathbb{R}^{2 \times (N^{\text{lo}}+N^{\text{up}}+2)}$ and the vector $\mathbf{d} \in \mathbb{R}^2$ are given by

$$\mathbf{C} := \begin{pmatrix} \mathbf{1}_{1,N^{\text{lo}}} & \mathbf{0}_{1,N^{\text{up}}} & 0 & 0 \\ \mathbf{0}_{1,N^{\text{lo}}} & \mathbf{1}_{1,N^{\text{up}}} & 0 & 0 \end{pmatrix}, \quad \mathbf{d} := \begin{pmatrix} 1 \\ 1 \end{pmatrix}.$$

By solving the linear program (7.37) we obtain an optimal vector \mathbf{x}_* , which defines a corresponding direction $\eta_{h,*} \in -\tilde{\Gamma}$ via (7.34) and (7.38). The element $\eta_{h,*}$ is expected to define a descent direction for the goal functional $\tilde{J} = J_{h,j} \circ Y_h \circ S_G$ at ξ_h . The choice strategy hence consists in choosing $-v_h := \eta_{h,*}$ for the update (7.24) of Algorithm 7.1.

We remark that linear programs of the form (7.37) can be solved e.g. by the MATLAB function `linprog`, which is provided by MATLAB's Optimization Toolbox. We also used the C++ function library `lp_solve`, which is available under a GNU Lesser General Public License (cf. lpsolve.sourceforge.net/5.5/).

7.7 Remarks

In Section 7.3 and Section 7.4 we discussed, that a photonic band structure optimization problems is, in fact, a constrained minimization problem. The so-called box constraint (7.14) arises from the definition of the extended admissible set $\bar{\mathcal{C}}_h$, and the symmetry constraint (7.21) arises from the assumption that approximate band structure computations can be restricted to quasimomentum vectors lying inside an irreducible zone of the first Brillouin zone. Algorithm 7.1 ensures both constraints by defining an alternative minimization, which is posed on the unconstrained space \mathcal{E}_h . In essence, this was accomplished through a parametrization of a subset of $\bar{\mathcal{C}}$ by virtue of the operator $Y_h \circ S_G$.

We also tested other methods for to ensure the optimization constraints. A simple way for ensuring the box constraint is to use a so-called *projective generalized gradient method*. Such a method was used e.g. by Cox and Dobson (cf. [29], [30]). A projective generalized gradient method constructs a minimizing sequence $\rho_h^{(0)}, \rho_h^{(1)}, \rho_h^{(2)}, \dots$, etc. in $\bar{\mathcal{C}}_h$ by choosing in each iteration a descent direction $-v_h^{(l-1)} \in -\Gamma_{J_h}(\xi^{(l-1)})$ and by applying the update rule

$$\rho_h^{(l)} := P_h(\rho_h^{(l-1)} - s_l v_h^{(l-1)}) \quad \text{for } l = 1, 2, \dots,$$

where the projection operator $P_h : \mathcal{E}_h \rightarrow \mathcal{C}_h$ is given by

$$P_h(\xi_h)(\mathbf{x}) := \begin{cases} \rho_{\min} & \text{if } \xi_h(\mathbf{x}) < \rho_{\min}, \\ \rho_{\max} & \text{if } \xi_h(\mathbf{x}) > \rho_{\max}, \\ \xi_h(\mathbf{x}) & \text{else} \end{cases} \quad \text{for all } \xi_h \in \mathcal{E}_h, \mathbf{x} \in \mathcal{X}_h.$$

We found that certain problems can occur for this approach: The choice strategies for the descent directions are usually based on a linear program similar to that in Section 7.6. Such linear programs, in essence, stem from a local linearization of the goal functional J_h . However, these linearizations cannot take into account the effect of the non-linear operator P_h . Therefore, it can happen that a descent direction for the functional J_h is chosen, which is not a descent direction for $J_h \circ P_h$. In such cases the optimization algorithm stops prematurely.

We also tested a so-called *penalty approach*, which consists in devising a strictly convex functional $g : \mathcal{E}_h \rightarrow [0, \infty)$, such that $g(\xi_h) = 0$ is satisfied if and only if $\xi_h \in \mathcal{C}$. Choosing a penalty parameter $a > 0$, one then considers the minimization problem for the functional $J_h + \alpha g$. We tested with different choices for α and g but could not obtain good results with this approach, so far.

Chapter 8

A Level-Set Method

In the previous chapter we developed an optimization algorithm for photonic band structure optimization problems. The algorithm was based on a so-called gradient method. We discussed that this method is not guaranteed to converge to an admissible solution. In this section we describe another algorithm, which is based on a so-called level set method. In Section 8.1 we introduce some preliminary concepts. In Section 8.1 we motivate an approach for a level set method for photonic band structure optimization problems. In Section 8.3 we describe a concrete optimization algorithm for such problems. In Section 8.4 we make some general remarks about this algorithm.

8.1 Basic Concepts

In this section we introduce some fundamental concepts of shape optimization, such as shape functionals, shape derivatives and level set methods. Here, we only give a brief overview of the most important concepts leaving out any mathematically rigorous derivations. For more details on shape optimization the reader can refer to the book of Sokołowski and Zolesio [70]. Level set methods are discussed in detail in the standard textbook of Osher and Fedkiw [58].

Suppose that Ω is a simply-connected, open domain in \mathbb{R}^n for some space dimension $n \in \mathbb{N}$, which is bounded by a Lipschitz smooth boundary curve $\partial\Omega$. In order to simplify notations, we define the linear function spaces

$$\mathcal{H} := H_{\text{per}}^1(\Omega, \mathbb{R}), \quad (8.1)$$

$$\mathcal{V} := W^{1,\infty}(\Omega, \mathbb{R}). \quad (8.2)$$

Here, $W^{1,\infty}(\Omega, \mathbb{R})$ denotes the Sobolev space that consist of all functions in $L^\infty(\Omega, \mathbb{R})$, which admit weak first partial derivatives in $L^\infty(\Omega, \mathbb{R})$.

Let $\mathcal{J} : \mathcal{H} \rightarrow \mathcal{P}(\Omega)$ be the mapping, which is densely defined by

$$\mathcal{J}(\varphi) := \{\mathbf{x} \in \Omega \mid \varphi(\mathbf{x}) < 0\} \quad \text{for all } \varphi \in C^\infty(\Omega, \mathbb{R}).$$

Here, $\mathcal{P}(\Omega)$ denotes the power set of Ω . We then define the system

$$\mathcal{S} := \{S \subset \Omega \mid \exists \varphi \in \mathcal{H} \text{ such that } S = \mathcal{J}(\varphi)\}. \quad (8.3)$$

A set $S \in \mathcal{S}$ is called an *admissible shape* in Ω . By definition, every admissible shape is an open subset of Ω . Furthermore, the boundary ∂S of each admissible shape $S \in \mathcal{S}$ coincides with the zero level set of a function in \mathcal{H} .

A function $f : \mathcal{S} \rightarrow \mathbb{R}$ is called a *shape functional*, and a minimization problem of the form

$$\underset{S \in \mathcal{S}}{\text{minimize}} \quad f(S) \quad (8.4)$$

is called a *shape optimization problem*. Problems of this type typically arise in structural optimization (see e.g. [3], [2]).

A *level set method* is an optimization method for shape optimization problems. The basic idea of a level set method is to represent an initial, admissible shape $S^{(0)} \in \mathcal{S}$ by a function $\varphi^{(0)} \in \mathcal{H}$ through

$$S^{(0)} = \mathcal{J}(\varphi^{(0)}).$$

By definition of \mathcal{S} , such a function $\varphi^{(0)}$ always exists. Now, suppose that $\psi \in C^1(\mathbb{R}, \mathcal{H})$ is a function that satisfies the first-order, parabolic initial boundary value problem

$$\begin{cases} \dot{\psi}(t, \cdot) + \mathbf{v}(t, \cdot) \cdot \nabla \psi(t, \cdot) = 0 & \text{in } \Omega \text{ for all } t \in (0, T], \\ \psi(0, \cdot) = \varphi^{(0)}. \end{cases} \quad (8.5)$$

for a given function $\mathbf{v} \in C^1(\mathbb{R}, \mathcal{V}^n)$ and a given positive number $T > 0$. Here, $\dot{\psi}(t)$ denotes the derivative of ψ with respect to the first variable, and $\nabla \psi(t, \cdot)$ denotes the gradient of the function $\psi(t, \cdot) \in \mathcal{H}$ for all $t \in \mathbb{R}$. Note that the initial boundary value problem is a diffusion-convection-type problem. Therefore, the function \mathbf{v} in (8.5) is commonly referred to as a *velocity field*. With the solution ψ of (8.5) at hand, one can define a family $\{S^{(t)}\}_{t \in [0, T]}$ of admissible shapes by

$$S^{(t)} := \mathcal{J}(\psi(t, \cdot)), \quad t \in [0, T]. \quad (8.6)$$

For every ‘‘point in time’’ $t \in [0, T]$ the function $\varphi^{(t)} := \psi(t, \cdot)$ is called the *level set function* of $S^{(t)}$. The initial boundary value problem (8.5) is commonly known as the *level set equation* for the family of level set functions $\{\varphi^{(t)}\}_{t \in [0, T]}$.

For a sufficiently small $T > 0$ the velocity field \mathbf{v} in (8.5) uniquely defines a vector field $\mathbf{T}_{\mathbf{v}} \in C^1([0, T], \mathcal{V}^n)$ as the solution of the initial value problem

$$\begin{cases} \mathbf{T}_{\mathbf{v}}(t, \mathbf{x}) = \mathbf{v}(t, \mathbf{T}_{\mathbf{v}}(t, \mathbf{x})) & \text{for all } t \in I, \mathbf{x} \in \Omega, \\ \mathbf{T}_{\mathbf{v}}(0, \mathbf{x}) = \mathbf{x} & \text{for all } \mathbf{x} \in \Omega, \end{cases}$$

such that $\mathbf{T}_{\mathbf{v}}(t, \cdot)$ is a $W^{1,\infty}$ -diffeomorphism for all $t \in [0, T]$ (see Theorem 2.16 in [70]). Furthermore, one can show (see Section 2.9 in [70]) that

$$S^{(t)} = \mathbf{T}_{\mathbf{v}}(t, S^{(0)}) \quad \text{for all } t \in [0, T],$$

where $S^{(t)}$ coincides with the admissible shape defined by (8.6).

Given a shape optimization problem of the form (8.4) and an initial admissible shape $S^{(0)} \in \mathcal{S}$, the basic idea of a level set method is to construct a velocity field $\mathbf{v} \in C^1(\mathbb{R}, \mathcal{V}^n)$, such that $f(S^{(t)}) = f(\mathbf{T}_{\mathbf{v}}(t, S^{(0)}))$ converges to a local minimum of f as t tends to ∞ . In order to construct such a velocity field, one needs to obtain some information about the local behaviour of $f(S^{(t)})$ for a given point in time $t \in \mathbb{R}$ with respect to small perturbations in t . Such information is provided by the concept of *shape differentiability*.

Given a shape functional $f : \mathcal{S} \rightarrow \mathbb{R}$ and a velocity field $\mathbf{v} \in C^1(\mathbb{R}, \mathcal{V}^n)$, the shape functional f is said to be *directionally shape differentiable* at a shape $S \in \mathcal{S}$ in direction \mathbf{v} , if the limit

$$f'(S; \mathbf{v}) := \lim_{t \rightarrow \infty} \frac{f(\mathbf{T}_{\mathbf{v}}(t, S)) - f(S)}{t}$$

exists. The limit $f'(S; \mathbf{v})$ is then called the *directional shape derivative* of f at S in direction \mathbf{v} . If $f'(S; \cdot)$ is a continuous, linear functional on $C^1(\mathbb{R}, \mathcal{V}^n)$, then f is called *shape differentiable* at S . The *shape differential*

$$Df(S) \in [C^1(\mathbb{R}, \mathcal{V}^n)]^*$$

of f at S is then defined by

$$Df(S)[\mathbf{v}] := f'(S; \mathbf{v}) \quad \text{for all } \mathbf{v} \in C^1(\mathbb{R}, \mathcal{V}^n).$$

We remark that $Df(S)[\mathbf{v}]$ is often referred to as the *shape derivative* of f at S .

8.2 Motivation

In this section we motivate a level set method for photonic band structure optimization problems, which involve the TE band structures of two-dimensional

photonic crystals. The motivation we present here is by no means mathematically rigorous. It is only meant to provide an intuitive understanding of how level set methods for photonic band structure optimization problems can be constructed.

In the following, we assume that Ω is a primitive domain of a Bravais lattice $\Lambda \subset \mathbb{R}^2$ of rank 2. By \mathbb{B} we denote the first Brillouin zone of Λ . Given an index $j \in \mathbb{N}$, a coefficient $\rho \in \mathcal{D}$ and a vector $\mathbf{k} \in \mathbb{B}$, the non-negative real number $\lambda_j^{\text{TE}}(\rho, \mathbf{k})$ was defined as the j -th smallest eigenvalue of the corresponding eigenvalue problem (4.94) on page 87. The function $u_j^{\text{TE}}(\rho, \mathbf{k}) \in V$ was defined as the corresponding eigenfunction. Without loss of generality, we assume that

$$m^{\text{TE}}(u_j^{\text{TE}}(\rho, \mathbf{k}), u_j^{\text{TE}}(\rho, \mathbf{k})) = 1 \quad \text{for all } j \in \mathbb{N}, \rho \in \mathcal{D}, \mathbf{k} \in \mathbb{B}.$$

Then, the j -th smallest TE eigenvalue $\lambda_j^{\text{TE}}(\rho, \mathbf{k})$ can be characterized as

$$\begin{aligned} \lambda_j^{\text{TE}}(\rho, \mathbf{k}) &= a_{\mathbf{k}}^{\text{TE}}(\rho)(u_j^{\text{TE}}(\rho, \mathbf{k}), u_j^{\text{TE}}(\rho, \mathbf{k})) \\ &= \int_{\Omega} \rho |(\nabla + i\mathbf{k})u_j^{\text{TE}}(\rho, \mathbf{k})|^2 \quad \text{for all } j \in \mathbb{N}, \rho \in \mathcal{D}, \mathbf{k} \in \mathbb{B}. \end{aligned}$$

For the purpose of this section we assume that every TE eigenvalue $\lambda_j^{\text{TE}}(\rho, \mathbf{k})$ is a geometrically simple eigenvalue, i.e., we assume that the eigenspace of every TE eigenvalue is one-dimensional.

We assume that the system \mathcal{S} of admissible shapes in Ω is defined according to (8.3). Then, we define the function set

$$\tilde{\mathcal{C}} := \{\rho_{\min} + (\rho_{\max} - \rho_{\min})\chi_S \mid S \in \mathcal{S}\}.$$

Clearly, $\tilde{\mathcal{C}}$ is a subset of the admissible set \mathcal{C} . Next, we define the mapping $\mathcal{R} : \mathcal{S} \rightarrow \tilde{\mathcal{C}}$ by

$$\mathcal{R}(S) := \rho_{\min} + (\rho_{\max} - \rho_{\min})\chi_S \quad \text{for all } S \in \mathcal{S}.$$

Given an index $j \in \mathbb{N}$ and a quasimomentum vector $\mathbf{k} \in \mathbb{B}$, we define the shape functional $\tilde{\lambda}_j^{\text{TE}}(\cdot, \mathbf{k}) : \mathcal{S} \rightarrow \mathbb{R}$ and the mapping $\tilde{u}_j^{\text{TE}}(\cdot, \mathbf{k}) : \mathcal{S} \rightarrow V$ by

$$\begin{aligned} \tilde{\lambda}_j^{\text{TE}}(\cdot, \mathbf{k}) &:= \lambda_j^{\text{TE}}(\cdot, \mathbf{k}) \circ \mathcal{R}, \\ \tilde{u}_j^{\text{TE}}(\cdot, \mathbf{k}) &:= u_j^{\text{TE}}(\cdot, \mathbf{k}) \circ \mathcal{R}. \end{aligned}$$

One easily verifies that the shape functional $\tilde{\lambda}_j^{\text{TE}}(\cdot, \mathbf{k})$ is given by

$$\begin{aligned} \tilde{\lambda}_j^{\text{TE}}(S, \mathbf{k}) &= \rho_{\min} \int_{\Omega} |(\nabla + i\mathbf{k})\tilde{u}_j^{\text{TE}}(S, \mathbf{k})|^2 \\ &\quad + (\rho_{\max} - \rho_{\min}) \int_S |(\nabla + i\mathbf{k})\tilde{u}_j^{\text{TE}}(S, \mathbf{k})|^2 \quad \text{for all } S \in \mathcal{S}. \end{aligned}$$

It can be shown that $\tilde{\lambda}_j^{\text{TE}}(\cdot, \mathbf{k})$ is shape differentiable at an admissible shape $S \in \mathcal{S}$, provided that $\lambda_j^{\text{TE}}(\tilde{\rho}, \mathbf{k})$ is a geometrically simple eigenvalue for all coefficients $\tilde{\rho}$ in a neighbourhood of $\mathcal{R}(S)$. The shape differential is given by

$$D\tilde{\lambda}_j^{\text{TE}}(\cdot, \mathbf{k})(S)[\mathbf{v}] = (\rho_{\max} - \rho_{\min}) \int_{\partial S} \boldsymbol{\nu}_S \cdot \mathbf{v}(0, \cdot) |(\nabla + i\mathbf{k})\tilde{u}_j^{\text{TE}}(S, \mathbf{k})|^2$$

for all $\mathbf{v} \in \mathcal{V}$, where $\boldsymbol{\nu}_S : \partial S \rightarrow \mathbb{S}^1$ denotes the outer unit normal field on the boundary of S .

In the following, we aim at constructing a minimizer $S_* \in \mathcal{S}$ for $\tilde{\lambda}_j^{\text{TE}}(\cdot, \mathbf{k})$. We want to construct the minimizer with a level set method. Therefore, we choose an arbitrary initial admissible shape $S_0 \in \mathcal{S}$, which is given by $S_0 = \mathcal{J}(\varphi^{(0)})$ for some initial function level set function $\varphi^{(0)} \in \mathcal{H}$. Given a velocity field $\mathbf{v} \in C^1(\mathbb{R}, \mathcal{V}^2)$, we can construct a family $\{S^{(t)}\}_{t \in [0, T]}$ of admissible shapes by solving the level set equation

$$\begin{cases} \dot{\psi}(t, \cdot) + \mathbf{v}(t) \cdot \nabla \psi(t, \cdot) = 0 & \text{in } \Omega \text{ for all } t \in (0, T], \\ \psi(0, \cdot) = \varphi^{(0)}. \end{cases} \quad (8.7)$$

for $\psi \in C^1([0, T], \mathcal{H})$, and by setting $S^{(t)} := \mathcal{J}(\psi(t, \cdot))$ for all $t \in [0, T]$. Clearly, this specific family of admissible shapes is completely determined by the initial level set function $\varphi^{(0)}$ and by the velocity field \mathbf{v} . Furthermore, one can show that the shape differential of $\tilde{\lambda}_j^{\text{TE}}(\cdot, \mathbf{k})$ satisfies

$$\begin{aligned} D\tilde{\lambda}_j^{\text{TE}}(\cdot, \mathbf{k})(S^{(t)})[\mathbf{v}] &= (\rho_{\max} - \rho_{\min}) \int_{\partial S^{(t)}} \boldsymbol{\nu}_{S^{(t)}} \cdot \mathbf{v}(t, \cdot) |(\nabla + i\mathbf{k})\tilde{u}_j^{\text{TE}}(S^{(t)}, \mathbf{k})|^2 \\ &\quad + o(t^2). \end{aligned} \quad (8.8)$$

Recall that every shape $S^{(t)}$ is determined by its corresponding level set function $\varphi^{(t)} = \psi(t, \cdot)$. Provided that the gradient $\nabla \varphi^{(t)}$ does not vanish on $\partial S^{(t)}$ for any $t \in [0, T]$, we can represent the outer unit normal fields $\boldsymbol{\nu}_{S^{(t)}} : \partial S^{(t)} \rightarrow \mathbb{S}^1$ on the boundaries of the admissible shapes $S^{(t)}$ by

$$\boldsymbol{\nu}_{S^{(t)}} = \frac{\nabla \varphi^{(t)}}{|\nabla \varphi^{(t)}|} \Big|_{\partial S^{(t)}} \quad \text{for all } t \in [0, T].$$

Under the further assumption that $\nabla \varphi^{(t)}$ does not vanish almost everywhere on Ω for any $t \in [0, T]$, we can define continuous extensions $\boldsymbol{\nu}^{(t)} : \Omega \rightarrow \mathbb{S}^1$ of the outer unit normal fields $\boldsymbol{\nu}_{S^{(t)}}$ by

$$\boldsymbol{\nu}^{(t)} = \frac{\nabla \varphi^{(t)}}{|\nabla \varphi^{(t)}|} \quad \text{for all } t \in [0, T].$$

Now, suppose that for every $t \in [0, T]$ the velocity field $\mathbf{v}(t, \cdot)$ is normal to $\partial S^{(t)}$, i.e., that $\mathbf{v}(t, \cdot)$ is given by

$$\mathbf{v}(t, \cdot) = c(t, \cdot) \boldsymbol{\nu}^{(t)} \quad \text{for all } t \in [0, T] \quad (8.9)$$

for some function $c \in C^1(\mathbb{R}, \mathcal{V})$. Then, (8.8) can be rewritten as

$$\begin{aligned} D\tilde{\lambda}_j^{\text{TE}}(\cdot, \mathbf{k})(S^{(t)})[\mathbf{v}] &= (\rho_{\max} - \rho_{\min}) \int_{\partial S^{(t)}} c(t, \cdot) |(\nabla + \mathbf{i}\mathbf{k})\tilde{u}_j^{\text{TE}}(S^{(t)}, \mathbf{k})|^2 \\ &\quad + o(t^2). \end{aligned} \quad (8.10)$$

Furthermore, the level set equation (8.7) reduces to a so-called *Hamilton–Jacobi*-type system, which reads

$$\begin{cases} \dot{\psi}(t, \cdot) + c(t, \cdot) |\nabla \psi(t, \cdot)| = 0 & \text{in } \Omega \text{ for all } t \in (0, T], \\ \psi(0, \cdot) = \varphi^{(0)}. \end{cases} \quad (8.11)$$

The function c in (8.11) is commonly referred to as the *speed field*.

Recall that, for every $t \in [0, T]$, the shape differential $D\tilde{\lambda}_j^{\text{TE}}(\cdot, \mathbf{k})(S^{(t)})$ is a linearization of the shape functional $\tilde{\lambda}_j^{\text{TE}}(\cdot, \mathbf{k})$ at the admissible shape $S^{(t)}$ with respect to perturbations of $S^{(t)}$ that are generated by velocity fields $\mathbf{v} \in C^1(\mathbb{R}, \mathcal{V}^2)$. Hence, a negative sign of the integral in (8.10) indicates, that the value of $\tilde{\lambda}_j^{\text{TE}}(S^{(t)}, \mathbf{k})$ is likely to decrease, if the shape $S^{(t)}$ is perturbed by the velocity \mathbf{v} , which is defined in (8.9). Clearly, a negative sign of the integral can be guaranteed by setting

$$c(t, \cdot) := -(\rho_{\max} - \rho_{\min}) |(\nabla + \mathbf{i}\mathbf{k})\tilde{u}_j^{\text{TE}}(S^{(t)}, \mathbf{k})|^2 \quad \text{for all } t \in [0, T].$$

Now, let $J^{\text{TE}} : \mathcal{D} \rightarrow \mathbb{R}$ be a given TE goal functional. We then define a corresponding shape functional $\tilde{J}^{\text{TE}} : \mathcal{S} \rightarrow \mathbb{R}$ by $\tilde{J}^{\text{TE}} := J^{\text{TE}} \circ \mathcal{R}$ and consider the minimization problem

$$\underset{S \in \mathcal{S}}{\text{minimize}} \quad \tilde{J}^{\text{TE}}(S).$$

Suppose that the goal functional J^{TE} is strictly differentiable at some density function $\rho = \mathcal{R}(S)$, where $S \in \mathcal{S}$ is an admissible shape, and that the strict differential of J^{TE} at ρ is given by

$$D_s J^{\text{TE}}(\rho)[\eta] = \int_{\Omega} \gamma \eta \quad \text{for all } \eta \in \mathcal{E} \quad (8.12)$$

for some uniquely defined function $\gamma \in L^1(\Omega)$. Then, one can show that the shape differential of \tilde{J}^{TE} at the admissible shape S is given by

$$D\tilde{J}^{\text{TE}}(S)[\mathbf{v}] = (\rho_{\max} - \rho_{\min}) \int_{\partial S} \boldsymbol{\nu}_S \cdot \mathbf{v}(0, \cdot) \gamma \quad \text{for all } \mathbf{v} \in C^1(\mathbb{R}, \mathcal{V}^2).$$

Starting off with an initial admissible shape $S^{(0)} \in \mathcal{S}$, one can attempt to construct a family $\{S^{(t)}\}_{t \in [0, T]}$ of admissible shapes, such that $\tilde{J}^{\text{TE}}(S^{(T)}) < J^{\text{TE}}(S^{(0)})$ for some $T > 0$. Using a level set method, this family is constructed by choosing an initial level set function $\varphi^{(0)}$, such that $S^{(0)} = \mathcal{J}(\varphi^{(0)})$. The level set function is then evolved by solving the Hamilton–Jacobi equation (8.11). According to the discussion above, it is reasonable to choose the speed field $c \in C^1(\mathbb{R}, \mathcal{V})$, such that

$$c(t, \cdot) = -(\rho_{\max} - \rho_{\min}) \gamma^{(t)} \quad \text{for all } t \in [0, T],$$

where $\gamma^{(t)} \in L^1(\Omega, \mathbb{R})$ is the function that determines the strict derivative of J^{TE} at $\rho^{(t)} := \mathcal{R}(S^{(t)})$ according to (8.12).

If the goal functional J^{TE} is not strictly differentiable but only locally Lipschitz continuous, it is necessary to adapt the construction principle for the level set method presented above. We hence construct the family $\{S^{(t)}\}_{t \in [0, T]}$ of admissible shapes as follows: Given a family member $S^{(t)}$ and a corresponding density function $\rho^{(t)} := \mathcal{R}(S^{(t)})$, we compute a convex set of functions $\Gamma(\rho^{(t)}) \subset L^1(\Omega, \mathbb{R})$, such that

$$\partial J^{\text{TE}}(\rho^{(t)}) \subset \left\{ p \in \mathcal{E}^* \mid p(\eta) = \int_{\Omega} \gamma \eta, \gamma \in \Gamma(\rho^{(t)}) \right\}.$$

The speed field $c \in C^1(\mathbb{R}, \mathcal{V})$ is then chosen, such that

$$c(t, \cdot) \in -(\rho_{\max} - \rho_{\min}) \Gamma(\rho^{(t)}) \quad \text{for all } t \in [0, T].$$

To the best of our knowledge there is no conclusive evidence that this approach is justified in all cases. Nevertheless, we were able to develop an optimization algorithm based on this approach, which performed quite well. We described this algorithm in the next section.

8.3 The Algorithm

Based on the assumptions in the previous section, we implemented a level set algorithm for photonic band structure optimization problems. In this section we describe this algorithm, which is also listed in pseudo-code as Algorithm 8.1.

Under the assumptions of Section 7.2, we further assume that the Hilbert space \mathcal{H} is discretized by a conforming finite element ansatz on the underlying mesh \mathcal{T}_h . We denote the corresponding finite element space by \mathcal{H}_h . We then define the so-called *interpretation operator* $\Theta_h : \mathcal{H} \rightarrow \mathcal{C}_h$ according to

$$\Theta_h(\varphi_h)(\mathbf{x}) := \begin{cases} \rho_{\min} & \text{if } \varphi_h(\mathbf{x}) > 0, \\ \rho_{\max} & \text{if } \varphi_h(\mathbf{x}) < 0, \end{cases} \quad \text{for all } \mathbf{x} \in \mathcal{X}_h.$$

Algorithm 8.1 Level-Set Method

Input:

$J_h : \mathcal{D}_h \rightarrow \mathbb{R}$	<i>Goal functional</i>
G	<i>Group of orthogonal transformations (Point group)</i>
$\varphi_h^{(0)} \in \mathcal{H}_h$	<i>Initial level set function</i>

Parameters:

$t_{\min} > 0$	<i>Minimal time step</i>
$t_{\max} > t_{\min}$	<i>Maximal step size</i>
$\Delta_{\min} > 0$	<i>Minimal decrease</i>

Algorithm:

```

 $\rho_h^{(0)} := \Theta_h(\varphi_h^{(0)})$ 
compute  $z_0 := J_h(\rho_h^{(0)})$ 
 $t_0 := t_{\max}$ 
for  $l = 1, 2, \dots$  do
  compute  $\Gamma_{J_h}(\rho_h^{(l-1)})$ 
  choose a descent direction  $-\gamma_h^{(l-1)} \in -(\rho_{\max} - \rho_{\min})\Gamma_{J_h}(\rho_h^{(l-1)})$ 
   $\gamma_h^{(l-1)} := \tilde{\gamma}_h^{(l-1)} / \|\tilde{\gamma}_h^{(l-1)}\|_{\infty}$ 
   $t_{l,1} := \min\{2t_{l-1}, t_{\max}\}$ 
  for  $p = 1, 2, \dots$  do
     $\varphi_h^{(l,p)} := \Psi_h(\varphi_h^{(l-1)}, -\gamma_h^{(l-1)}, t_{l,p})$ 
     $\rho_h^{(l,p)} := \Theta_h(\varphi_h^{(l,p)})$ 
    compute  $z_{l,p} := J_h(\rho_h^{(l,p)})$ 
    if  $z_{l-1} - z_{l,p} > \Delta_{\min}$  then
       $\varphi_h^{(l)} := \varphi_h^{(l,p)}$ 
       $\rho_h^{(l)} := \rho_h^{(l,p)}$ 
       $z_l := z_{l,p}$ 
       $t_l := t_{l,p}$ 
      break
    end if
     $t_{l,p+1} := t_{l,p}/2$ 
    if  $t_{l,p+1} < t_{\min}$  then
      return  $\rho_h^{(l-1)}$ 
    end if
  end for
end for
end for

```

Here, the understanding is that the function φ_h represents the level set function of a shape $S \in \mathcal{S}$. The shape S itself is understood to represent the area inside the primitive domain of a photonic crystal with optical density ρ_{\max} .

For convenience, we define the so-called *evolution operator* $\Psi : \mathcal{H} \times \mathcal{V} \times \mathbb{R} \rightarrow \mathcal{H}$ as follows. Given a level set function $\varphi \in \mathcal{H}$, a speed field $c \in \mathcal{V}$, and a positive number $T > 0$, let $\psi \in C^1([0, T], \mathcal{H})$ be the solution of the initial boundary value problem

$$\begin{cases} \dot{\psi}(t, \cdot) + c|\nabla\psi(t, \cdot)| = 0 & \text{in } \Omega \text{ for all } t \in (0, T], \\ \psi(0, \cdot) = \varphi. \end{cases}$$

Then, we define $\Psi(\varphi, c, T) := \psi(t, \cdot)$. In practice, the evolution operator Ψ is realized by a *discretized evolution operator* $\Psi_h : \mathcal{H}_h \times \mathcal{E}_h \times \mathbb{R} \rightarrow \mathcal{H}_h$. Such a discretized evolution operator is evaluated by solving the above initial boundary value problem numerically. In two space dimensions we implemented a second-order, semi-discrete central evolution scheme for Hamilton-Jacobi equations, which was proposed by Kurganov and Tadmor (see [51]).

Given a discretized band structure optimization problem of the form

$$\underset{\rho_h \in \mathcal{C}_h}{\text{minimize}} \quad J_h(\rho_h),$$

we assume that the initial density function $\rho_h^{(0)}$ is determined by an *initial level set function* $\varphi_h^{(0)}$, which we require as an input parameter for the algorithm. In every iteration, a superset $\Gamma_{J_h}(\rho_h^{(l-1)})$ of the generalized gradient of J_h at $\rho_h^{(l-1)}$ is computed. A descent direction $-\gamma_h^{(l-1)}$ is then chosen from the set $-\Gamma_{J_h}(\rho_h^{(l-1)})$. The choice is accomplished through a suitable choice strategy. We found that a choice strategy similar to that developed in Section 7.6 is suitable, if the goal functional J_h is given by a gap width functional as defined in Section 5.5. A new discretized level set function $\varphi_h^{(l)}$ is obtained as the result of $\Psi_h(\varphi_h^{(l-1)}, -\gamma_h^{(l-1)}, t_{l,p})$.

The “time step” $t_{l,p}$ of the level set evolution is determined adaptively by a trial-and-error approach. Whenever an evolution by a given time step fails to produce a level set function that yields a prescribed minimal decrease Δ_{\min} of the goal value, the step size is divided by two. In order to avoid excessively small step sizes, the step size is doubled at the beginning of each iteration, without exceeding a prescribed maximal time step t_{\min} , however. If the time step drops below a prescribed minimal time step t_{\min} , the algorithm terminates. In this case the discretized level set function of the previous time step is chosen as the *final level set function* $\phi_h^{(L)}$. Through $\rho_h^{(L)} := \Theta_h(\varphi_h^{(L)})$ a final density function $\rho_h^{(L)}$ is determined, which is then returned by the algorithm. By construction, this final density function belongs to the admissible set \mathcal{C}_h .

We conducted a series of numerical experiments on photonic band gap maximization problems for two-dimensional photonic crystals using Algorithm 8.1. Some of the results are presented in Chapter 9.

8.4 Remarks

From an algorithmic point of view the generalized gradient method as given by Algorithm 7.1 and the level set method as given by Algorithm 8.1 are very similar. In both methods descent directions are chosen from a superset of the goal functional's generalized gradient. The methods only differ in how they use this descent direction to construct new density functions. In the generalized gradient method the descent directions are simply added to the current density functions. In contrast to this, the level method uses the descent direction as a speed function in a Hamilton–Jacobi equation by which the level set is evolved.

In the literature, it is often discussed that the level set functions can develop slopes, which are either very steep or very flat. Both cases are known to cause problems for the convergence of a level set method. It is often proposed to restart the level set method periodically. Another method consists in regularizing the level set function after several iterations by solving the problem

$$\begin{aligned} \dot{\psi}(t, \cdot) + \operatorname{sgn}(\varphi)(|\nabla\psi(t, \cdot)| - 1) &= 0 \quad \text{in } \Omega \text{ for all } t > 0, \\ \psi(0, \cdot) &= \varphi. \end{aligned}$$

This problem is known to converge to a so-called signed distance function for the zero level set of ϕ (see e.g. Section 7.4 in [58]). In our numerical experiments we did not implement either of the above techniques as convergence was not found to be a problem.

Chapter 9

Numerical Results

In this chapter we demonstrate the performance of the optimization algorithms developed in Chapter 7 and Chapter 8. In Section 9.1 and Section 9.2 we present numerical experiments on photonic band gap maximization problems (PBGMPs) for two-dimensional photonic crystals. Numerical experiment on PBGMPs for three-dimensional photonic crystals are presented in Section 9.3.

9.1 Maximizing TM Band Gaps

In this section we present some numerical solutions of photonic band gap maximization problems (PBGMPs) involving TE band gaps of two-dimensional photonic crystals. Each problem was given by a minimization problem of the form

$$\underset{\rho \in \mathcal{C}}{\text{minimize}} \quad J_j^{\text{TM}}(\rho),$$

where J_j^{TM} denotes the TM band gap functional defined by (5.14) for a given index $j \in \mathbb{N}$. Throughout this section, we only consider photonic crystals, whose medium structures are periodic with respect to a square lattice $a\mathbb{Z}^2$, where $a > 0$ denotes the lattice constant (see Section 3.2). Furthermore, we assume that the point group of each crystal is isomorphic to the orthogonal group $O_2(\mathbb{Z})$. In all numerical experiments the computational domain Ω was given by the square $(-a/2, a/2)^2$.

The photonic crystals that were considered consisted of exactly two different materials with relative electric permittivities $\varepsilon_{r,\min} = 1$ and $\varepsilon_{r,\max} = 9$. Hence, the admissible set was given by

$$\mathcal{C} := \{\rho \in \mathcal{E} \mid \rho = 1 + 8\chi_S, S \in \mathcal{M}(\Omega)\}.$$

We discretized the computational domain Ω by a 100×100 , rectangular, structured mesh. The approximate band structures were computed through an H^1 -conforming, bilinear finite element ansatz. The density functions were discretized

quadrature-point-wise as described in Section 7.2. We used a custom-built MATLAB toolbox to implement the finite element ansatz. This toolbox is briefly described in the Appendix Chapter A. The toolbox uses Lobatto-type shape functions (cf. in [73]) and Gauss-type quadrature rules. We used a rule of consistency order 7, which gave 12 quadrature points on each mesh element (cf. Table 4.40 in [73]). The finite element mesh as well as the quadrature rule were chosen such that such that the set of quadrature points \mathcal{X}_h was $O_2(\mathbb{Z})$ -symmetric, which facilitated the implementation of symmetrization operations. All eigenvalue problems were solved using MATLAB's `eigs` function.

9.1.1 Results of the Generalized Gradient Method

In the following examples we applied the generalized gradient method as given by Algorithm 7.1 to various PBGMPs. The parameters of Algorithm 7.1 were chosen as $s_{\min} = 10^{-6}$, $s_{\max} = 1$, and $\Delta_{\min} = 10^{-8}$. The funnel operator Y_h was given as defined by (7.15) and (7.20) with $r = 4$. When applied to an admissible density function in $\rho_h \in \mathcal{C}_h$, this funnel operator returns a two-valued function with function values of approximately 1.1439 and 7.8679. The descent directions were chosen by solving linear programs, as explained in Section 7.6. The linear programs were solved using the MATLAB function `linprog`, which is provided by the MATLAB Optimization Toolbox.

Since the generalized gradient method is not guaranteed to converge to an admissible, two-valued solution, we always comment on whether or not a solution can be considered to be close to an admissible solution. To make this notion precise, define for every positive number $\delta > 0$ the set

$$\Xi(\delta) := [\varepsilon_{r,\min}, \varepsilon_{r,\min} + \delta) \cup (\varepsilon_{r,\max} - \delta, \varepsilon_{r,\max}].$$

With this, a discretized permittivity function $\varepsilon_h \in \mathcal{E}_h$ with function values in the interval $[\varepsilon_{r,\min}, \varepsilon_{r,\max}]$ is considered to be close to an admissible permittivity function, if there exists a small number $\delta > 0$, such that ε_h takes function values in $\Xi(\delta)$ at most quadrature points. Recall that the discrete function space \mathcal{E}_h was defined by (7.5) on page 131. In order to obtain admissible permittivity functions, we subjected each solution of Algorithm 7.1 to the so-called *threshold operator* \tilde{H}_h , which is defined by

$$\tilde{H}_h \varepsilon_h(\mathbf{x}) := \begin{cases} \varepsilon_{r,\min} & \text{if } \rho(\mathbf{x}) < \frac{\varepsilon_{r,\min} + \varepsilon_{r,\max}}{2}, \\ \varepsilon_{r,\max} & \text{else .} \end{cases} \quad \text{for all } \varepsilon_h \in \mathcal{E}_h, \mathbf{x} \in \mathcal{X}_h, \quad (9.1)$$

We also remind the reader that Π_h was defined as the interpolation operator, which projects suitable functions onto the discrete function space \mathcal{E}_h (see Section 7.2).

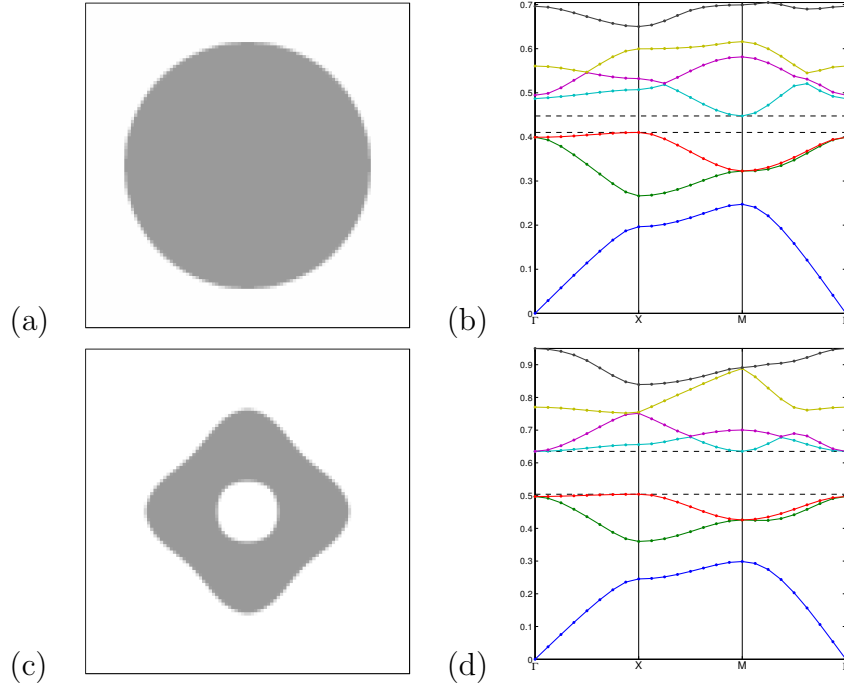


Figure 9.1: Results of Example 9.1. (a) Initial permittivity. (b) Initial band structure. (c) Final permittivity after 117 iterations. (d) Final band structure.

Example 9.1. In our first example we attempted to maximize a band gap in the TM band structure of a two-dimensional photonic crystal. The crystal consisted of a periodic arrangement of cylindrical rods. The radius of each rod was given by $0.38a$, where a denotes the lattice constant. The relative electric permittivity of the rods was given by $\varepsilon_{r,\max} = 9$, that of the ambient medium by $\varepsilon_{r,\min} = 1$. This medium structure is represented by the relative electric permittivity function ε_r , whose restriction to Ω is given by

$$\varepsilon_r|_{\Omega} = 1 + 8\chi_{B_{0.38a}(\mathbf{0})},$$

where $B_r(\mathbf{x})$ denotes the open ball in \mathbb{R}^2 with radius $r > 0$ and center point $\mathbf{x} \in \mathbb{R}^2$. The TM band structure of the photonic crystal exhibited a band gap between the third and the fourth band. We computed an approximate gap width of $w_{3,h} = 0.0444 \cdot 2\pi c_0/a$, where c_0 denotes the speed of light.

Choosing $\xi_h^{(0)} := \mathcal{I}_h \varepsilon_r|_{\Omega}$ as the initial free density function, we obtained an initial permittivity function $\varepsilon_{r,h}^{(0)}$, which yielded a smaller gap width of approximately $w_{3,h}^{(0)} = 0.0374 \cdot 2\pi c_0/a$ (see Figure 9.1(a) and (b)). After 117 iterations, however, the algorithm converged to a final permittivity function $\varepsilon_{r,h}^{(117)}$,

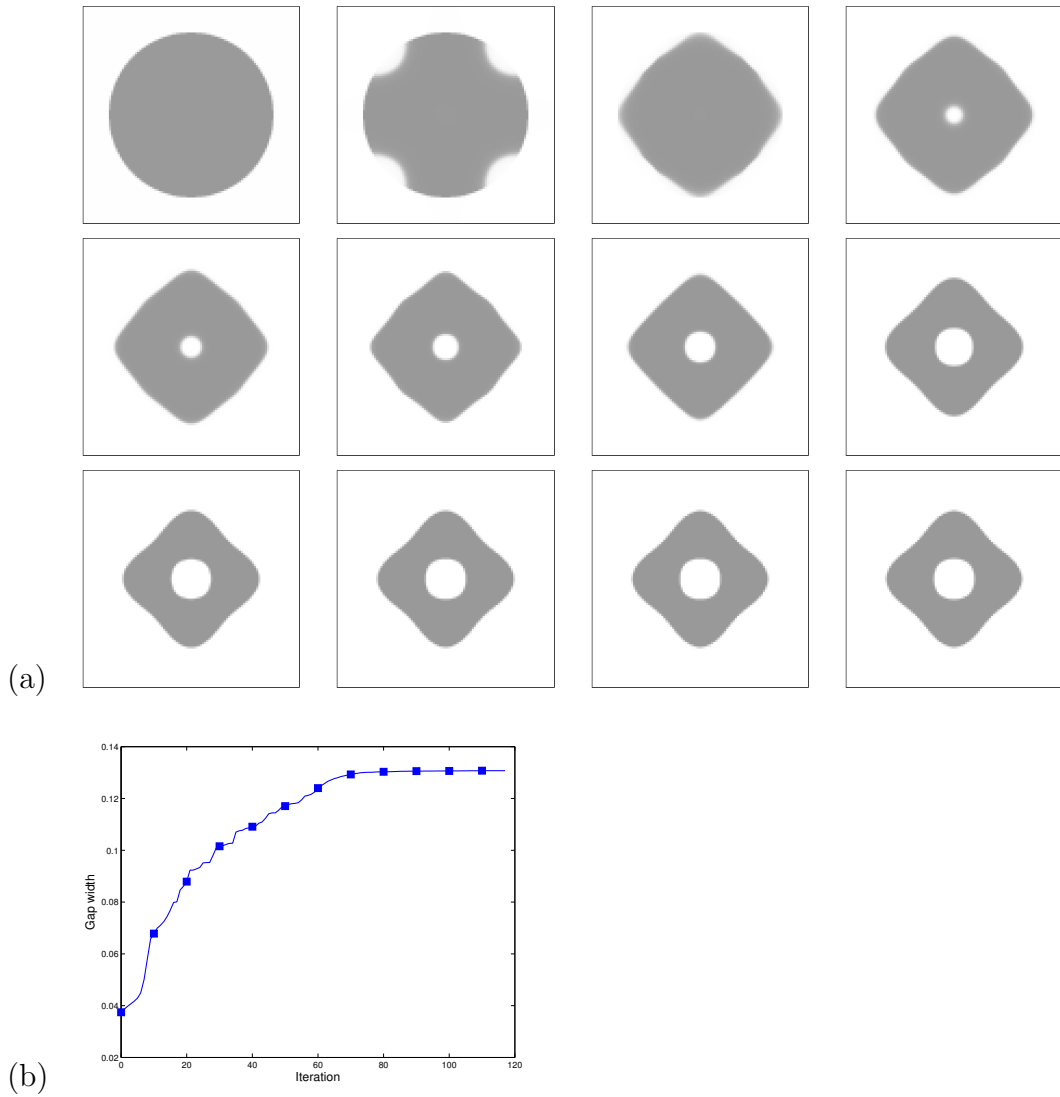


Figure 9.2: The evolution of the discretized permittivity function (a) and the corresponding gap width (b) in Example 9.1.

which represented a periodic arrangement of hollow rods with flower-shaped cross-sections (see Figure 9.1(c)). The width of the band gap increased to $w_{3,h}^{(117)} = 0.1307 \cdot 2\pi c_0/a$ (see Figure 9.1(d)).

We found that the minimal and maximal function values of $\varepsilon_{r,h}^{(117)}$ were given by 1.0000 and 8.9998, respectively. Moreover, $\varepsilon_{r,h}^{(117)}$ took function values in $\Xi(10^{-4})$ at 54.% of all quadrature points. Function values in $\Xi(10^{-3})$ and in $\Xi(10^{-2})$ were attained by $\varepsilon_{r,h}^{(117)}$ at 93.5% and 99.3% of all quadrature points,

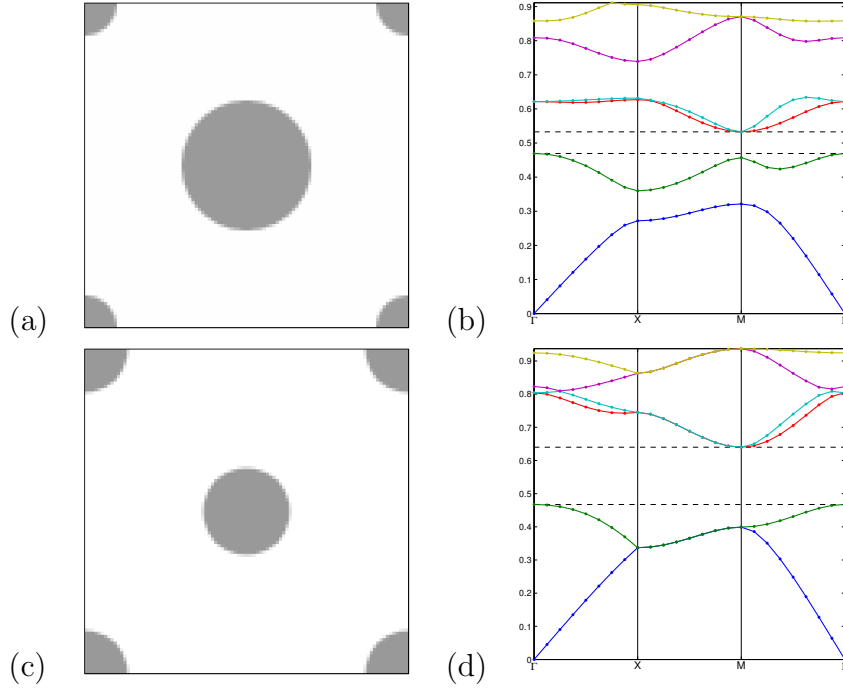


Figure 9.3: Results of Example 9.2. (a) Initial permittivity. (b) Initial band structure. (c) Final permittivity after 133 iterations. (d) Final band structure.

respectively. The results indicate that the final permittivity function is reasonably close to an admissible, two-valued function. By applying the threshold operator \tilde{H}_h to $\varepsilon_{r,h}^{(117)}$, we obtained such a function $\tilde{\varepsilon}_{r,h}$, which yielded a slightly larger gap width of $\tilde{w}_{3,h} = 0.1308 \cdot 2\pi c_0/a$.

Figure 9.2(a) shows the evolution of the permittivity function. From left to right and from top to bottom the images in Figure 9.2(a) show the discretized permittivity functions after every tenth iteration, starting with the initial permittivity function. Figure 9.2(b) shows the corresponding evolution of the gap width. The solid squares indicate points on the graph, which correspond to images in Figure 9.2(a).

We remark that the results in this example correspond well to results published by Cox and Dobson (cf. Figure 5.1 in [29]). Using a different generalized gradient method and a different goal functional, Cox and Dobson found a final permittivity function similar to ours, which yielded an approximate gap width of $0.128 \cdot 2\pi c_0/a$. In contrast to our generalized gradient method, the algorithm by Cox and Dobson needed a total of 1620 iterations to converge.

Example 9.2. In this example we maximize the band gap between the second and the third TM band of a two-dimensional photonic crystal. The crystal we start off with is composed of a staggered arrangement of cylindrical rods with two different radii. Each rod with a larger radius of $0.2a$ is surrounded by four rods with a smaller radius of $0.1a$. This medium structure is represented by a relative electric permittivity function, which satisfies

$$\varepsilon_r|_{\Omega} = 1 + 8\chi_{B_{0.2a}(\mathbf{0})} + \sum_{x_1, x_2 \in \{-a/2, a/2\}} \chi_{B_{0.1a}(x_1\mathbf{e}^{(1)} + x_2\mathbf{e}^{(2)}) \cap \Omega},$$

where $\mathbf{e}^{(1)}$ and $\mathbf{e}^{(2)}$ denote the standard basis vectors in \mathbb{R}^2 . The width of the band gap between the second and third TM band is approximately equal to $w_{2,h} = 0.0732 \cdot 2\pi c_0/a$.

As in the previous example, we chose $\mathcal{I}_h \varepsilon_r|_{\Omega}$ as the initial free density $\xi^{(0)}$. The corresponding initial gap width was given by $w_{2,h}^{(0)} = 0.0634 \cdot 2\pi c_0/a$ (see Figure 9.3(a) and (b)). After 133 iterations Algorithm 7.1 converged to a final permittivity function $\varepsilon_{r,h}^{(133)}$ representing a two-dimensional photonic crystal consisting of a periodic arrangement of identical cylindrical rods. The radius of each rod was approximately equal to $0.125a$. The gap width increased to $w_{2,h}^{(133)} = 0.1725 \cdot 2\pi c_0/a$ (see Figure 9.3(c) and (d)). This result corresponds well to a result published by Koa, Osher, and Santosa (cf. Figure 3 in [59]), who used a level set method and a different value for $\varepsilon_{r,\max}$.

The minimal and maximal function values of the final relative permittivity function $\varepsilon_{r,h}^{(133)}$ were 1.0000 and 9.0000. Remarkably, the function $\varepsilon_{r,h}^{(133)}$ was virtually two-valued in the sense that it took function values in $\Xi(10^{-4})$ at 99.75% of all quadrature points. We applied the threshold operator \tilde{H}_h to $\varepsilon_{r,h}^{(133)}$ and found that the resulting two-valued permittivity function yielded a gap width $\tilde{w}_{2,h}$, which satisfied $|\tilde{w}_{2,h} - w_{2,h}^{(133)}|/\tilde{w}_{2,h} < 10^{-6}$.

9.1.2 Results of the Level Set Method

Here, we present results obtained by the level set method as given by Algorithm 8.1. The parameters for this algorithm were chosen as $t_{\min} = 0.01$, $t_{\max} = 0.1$, $\Delta_{\min} = 10^{-8}$. Descent direction were chosen by a choice strategy similar to that described in Section 7.6.

Example 9.3. In this example we maximize the band gap between the first and the second TM band of a two-dimensional photonic crystal, which is composed of periodically placed, square rods. The cross-section of each square rod has an side length of $0.6a$. The TM band structure of this photonic crystal exhibits a band

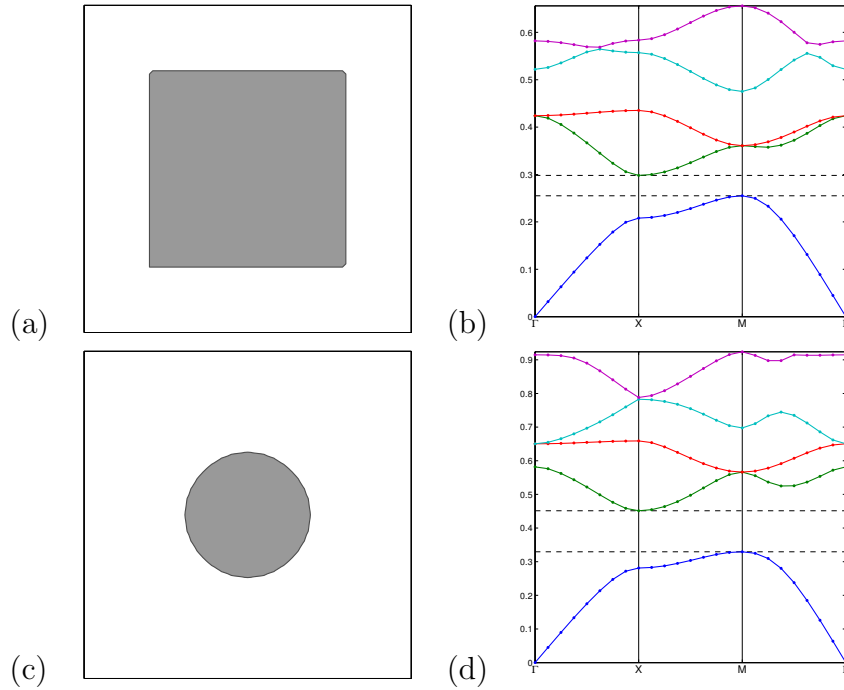


Figure 9.4: Results of Example 9.3. (a) Initial permittivity. (b) Initial band structure. (c) Final permittivity after 117 iterations. (d) Final band structure.

gap with a width of approximately $w_{1,h} = 0.0145 \cdot 2\pi c_0/a$ between the first and the second band (see Figure 9.4(a) and (b)).

We represented the crystal's medium structure in the primitive domain Ω by an initial level set function $\varphi^{(0)}$, which was given by

$$\varphi^{(0)}(\mathbf{x}) := 0.3a - |\mathbf{x}|_\infty \quad \text{for all } \mathbf{x} \in \Omega.$$

Here, $|\cdot|_\infty$ denotes the supremum norm in \mathbb{R}^2 . We remark that $\varphi^{(0)}$ is a signed distance function for the boundary of the rod's cross-section. After only 14 iterations Algorithm 8.1 converged to a final level set function $\varphi^{(14)}$, which represented the medium structure of a photonic crystal consisting of a periodic arrangement of cylindrical rods. The radius of each rod was approximately equal to $0.19a$. We computed an approximate gap width of $w_{1,h}^{(14)} = 0.1219 \cdot 2\pi c_0/a$ for this final medium structure (see Figure 9.4(c) and (d)). A similar optimal medium structure was also found by Kao, Osher, and Yablonovitch, though for a different relative electric permittivity of the rods (cf. Figure 2 in [44]).

In Figure 9.5(a) we depicted the permittivity functions after every second iteration of Algorithm 8.1. Figure 9.5 shows the evolution of the gap width. Each solid square corresponds to a permittivity function in Figure 9.5(a). Notice how

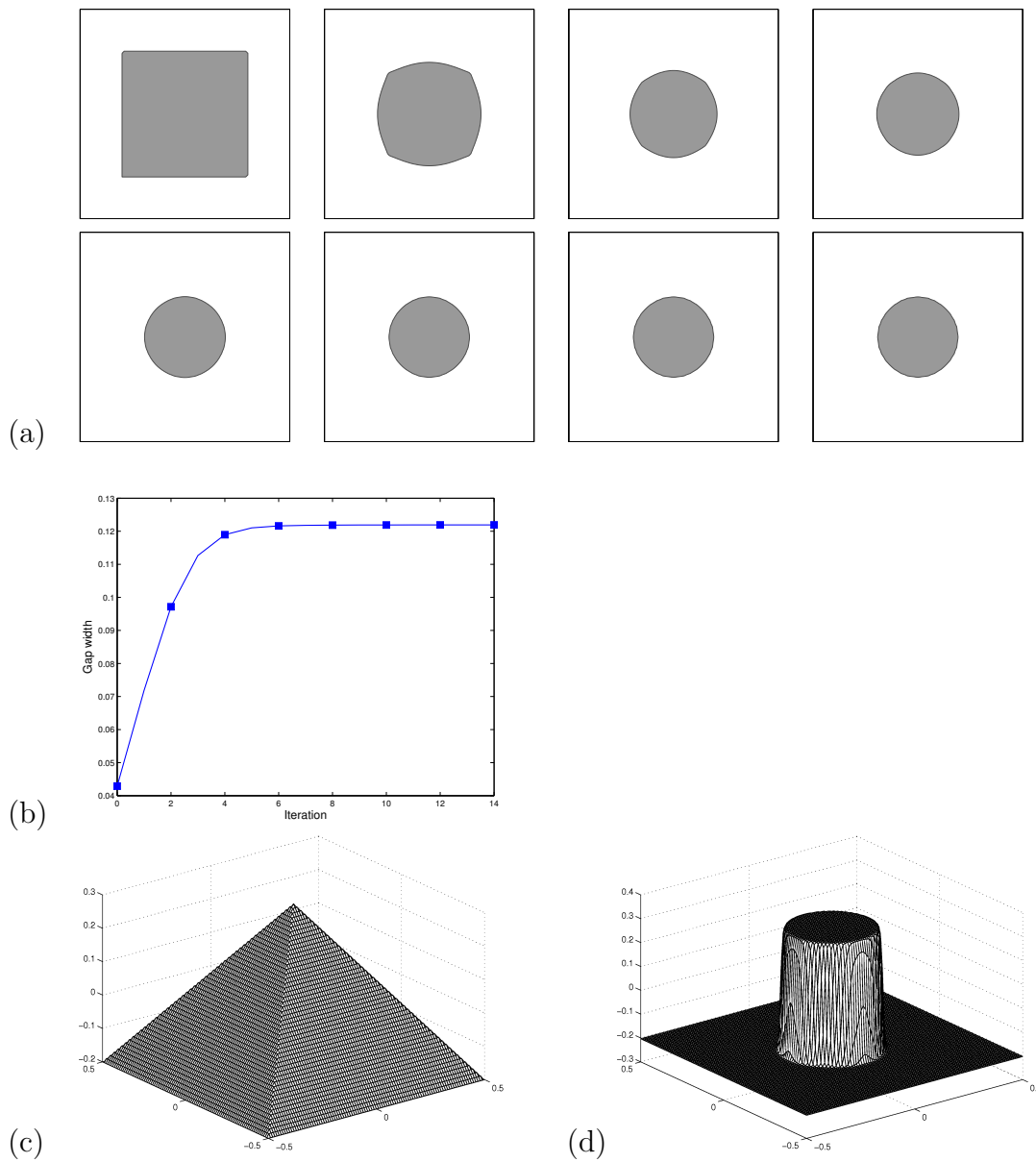


Figure 9.5: The evolution of the permittivity function (a) and the corresponding gap width (b) in Example 9.3.

fast the level set algorithm changes the rod's cross-section to a circle-like shape in the first few iterations. We compared the performance of Algorithm 8.1 to that of Algorithm 7.1, the generalized gradient method, and found that the latter needed 141 iterations to converge to the same final permittivity function.

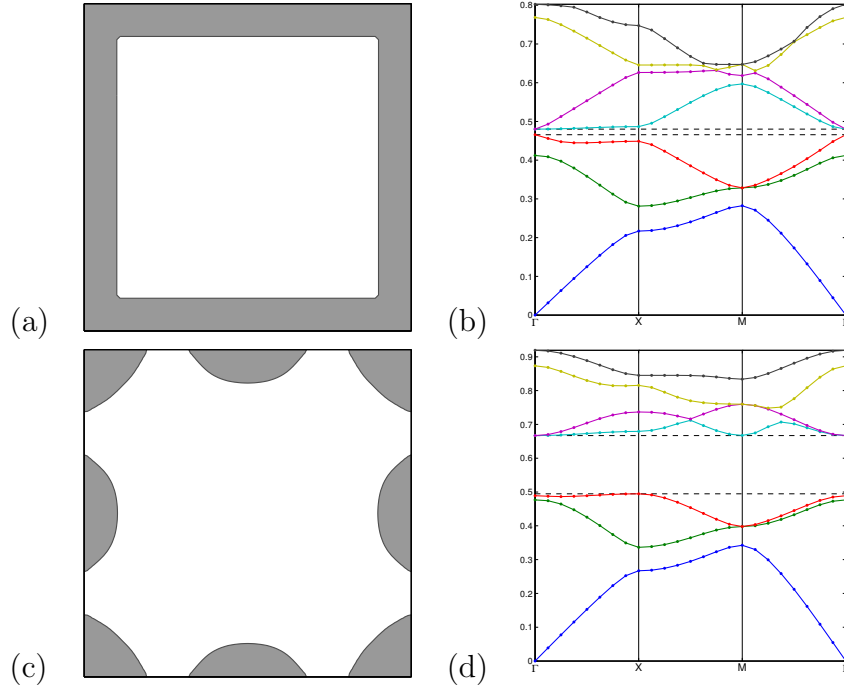


Figure 9.6: Results of Example 9.4. (a) Initial permittivity. (b) Initial band structure. (c) Final permittivity after 6 iterations. (d) Final band structure.

Example 9.4. In this example we start off with a photonic crystal, whose medium structure is given by a square lattice framework. The width of the framework is given by $0.2a$. The crystal's TM band structure exhibits a band gap between the third and the fourth band with an approximate gap width of $w_{3,h} = 0.0145 \cdot 2\pi c_0/a$ (see Figure 9.6(a) and (b)). We represented the medium structure by an initial level set function $\varphi^{(0)}$, given by

$$\varphi^{(0)}(\mathbf{x}) := |\mathbf{x}|_\infty - 0.4a \quad \text{for all } \mathbf{x} \in \mathbb{R}^2.$$

Only 6 iterations were needed to widen the band gap to an approximate gap width of $w_{3,h}^{(6)} = 0.1723 \cdot 2\pi c_0/a$ (see Figure 9.6(c) and (d)). This gap width is even larger than the one found in Example 9.1. We remark that our final medium structure corresponds well to optimal structures found by Kao, Osher, and Yablonovitch (cf. Figure 4 in [44]), as well as by Cox and Dobson (cf. Figure 3 in [29]).

Example 9.5. In this example we demonstrate that the level set method given by Algorithm 8.1 can also be used to open band gaps between separated photonic

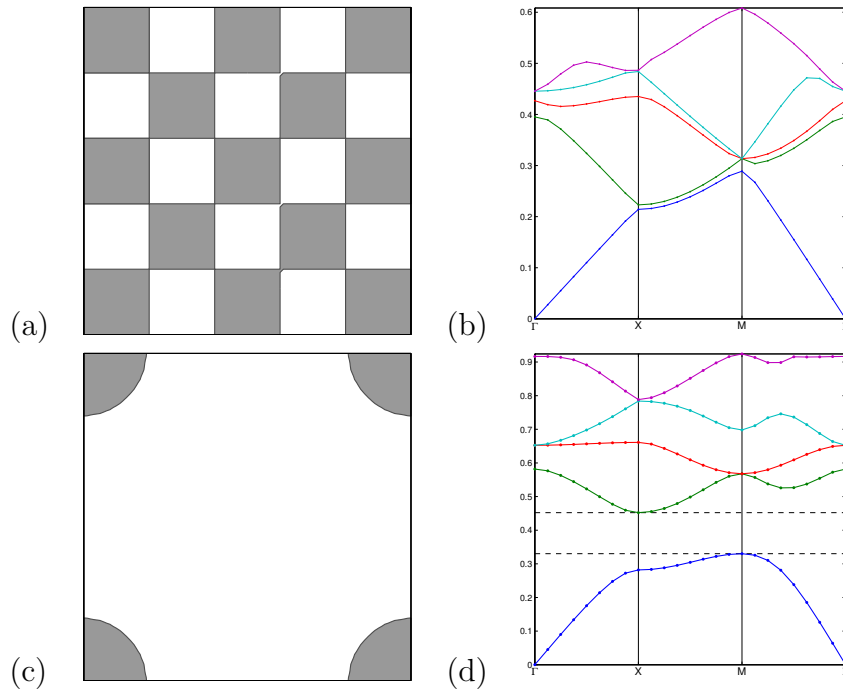


Figure 9.7: Results of Example 9.5. (a) Initial permittivity. (b) Initial band structure. (c) Final permittivity after 17 iterations. (d) Final band structure.

bands. By separated we mean that the bands do not intersect in any point. We considered a two-dimensional photonic crystal with a checkerboard-like medium structure. We chose the computational domain Ω , such that it contained a 5×5 portion of the checkerboard pattern (see Figure 9.7(a)). The TM band structure of the photonic crystal does not exhibit any band gap between the first 5 bands. One notices, however, that the first and the second band are separated (see Figure 9.7). We therefore applied Algorithm 8.1 with the gap width functional J_1^{TM} . We represented the initial medium structure by a level set function $\varphi^{(0)}$, which was given by

$$\varphi^{(0)}(\mathbf{x}) := \cos\left(\frac{5\pi\mathbf{x}_1}{a}\right) \cos\left(\frac{5\pi\mathbf{x}_2}{a}\right) \quad \text{for all } \mathbf{x} \in \mathbb{R}^2.$$

In the course of 17 iterations Algorithm 8.1 opened up a band gap between the first and the second TM band. The final level set function $\varepsilon^{(18)}$ represented the same optimized medium structure as in Example 9.3. It is therefore not surprising the final width of the band gap was approximately equal to $w_{1,h}^{(17)} = 0.1219 \cdot 2\pi c_0/a$.

In Figure 9.8 we depicted the change of the crystal's permittivity function during the first 7 iterations of the level set method. As one can see, the level set method had no trouble changing the medium structure's topology. The final

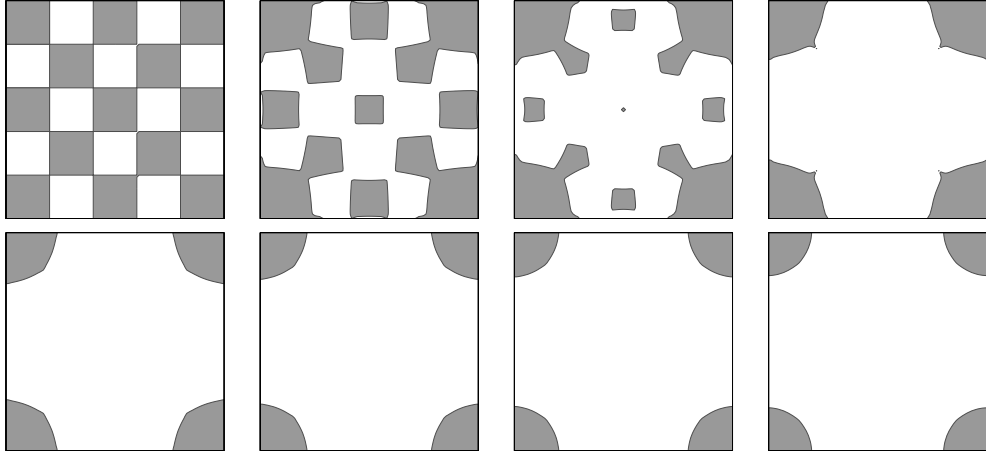


Figure 9.8: The evolution of the permittivity over the first 7 iterations in Example 9.5.

topology is already attained after the fourth iteration.

9.2 Maximizing TE Band Gaps

Here, we present numerical solutions of photonic band gap maximization problems involving the TE band structures of two-dimensional photonic crystals. Each such problem is given by a minimization problem

$$\underset{\rho \in \mathcal{C}^{\text{TE}}}{\text{minimize}} \quad J_j^{\text{TE}}(\rho),$$

where J_j^{TE} denotes the TM band gap functional defined by (5.14) for a given index $j \in \mathbb{N}$. In the following experiments the same type of photonic crystals were considered. With the two permittivities $\varepsilon_{r,\min} = 1$ and $\varepsilon_{r,\max} = 9$, the admissible is given by

$$\mathcal{C}^{\text{TE}} := \left\{ \rho \in \mathcal{E} \mid \rho = \frac{1}{9} + \frac{8}{9}\chi_S, \quad S \in \mathcal{M}(\Omega) \right\}.$$

The problems were discretized as described in the previous section.

9.2.1 Results of the Generalized Gradient Method

The results in the following examples were obtained by the Algorithm 7.1. The algorithm parameters were chosen as $s_{\min} = 10^{-6}$, $s_{\max} = 10^{-1}$, and $\Delta_{\min} = 10^{-8}$.

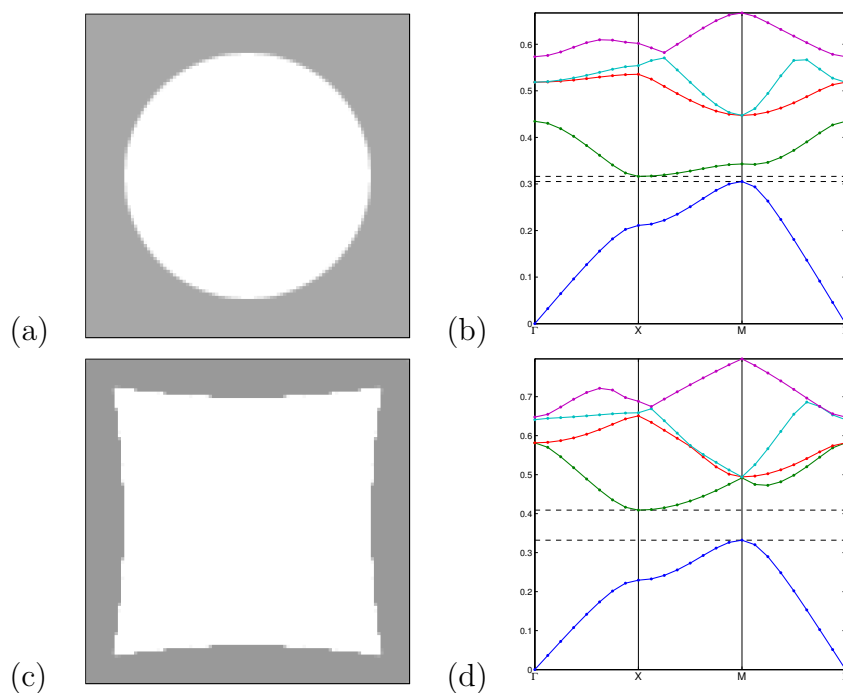


Figure 9.9: Results of Example 9.6. (a) Initial permittivity. (b) Initial band structure. (c) Final permittivity after 58 iterations. (d) Final band structure.

The funnel operator Y_h was given as defined by (7.15) and (7.20) with $r = 4$. When applied to an admissible density function in \mathcal{C}^{TE} , this funnel operator yields a two-valued density function, which represents a two-value permittivity function with function values of approximately 7.8679 and 1.0162. The descent directions were chosen as described in Section 7.6.

Example 9.6. In this example we maximize the band gap between the first and second TE band of a two-dimensional photonic crystal consisting of a periodic arrangement of cylindrical holes drilled into a material with relative electric permittivity $\varepsilon_{r,\max} = 9$. The medium structure is represented by the relative electric permittivity function ε_r , which satisfies

$$\varepsilon_r|_{\Omega} = 9 - 8B_{0.38a}(\mathbf{0}),$$

where a denotes the lattice constant, and we computed an approximate gap width of $w_{1,h} = 0.0158 \cdot 2\pi c_0/a$ for this medium structure.

Choosing $\xi_h^{(0)} := \mathcal{I}_h(1/\varepsilon_r|_{\Omega})$, the approximate initial gap width was given by $w_{1,h}^{(0)} = 0.0110 \cdot 2\pi c_0/a$ (see Figure 9.9(a) and (b)). Over the course of 58 iterations the gap width increased to a value of $w_{1,h}^{(58)} = 0.0771 \cdot 2\pi c_0/a$. The final

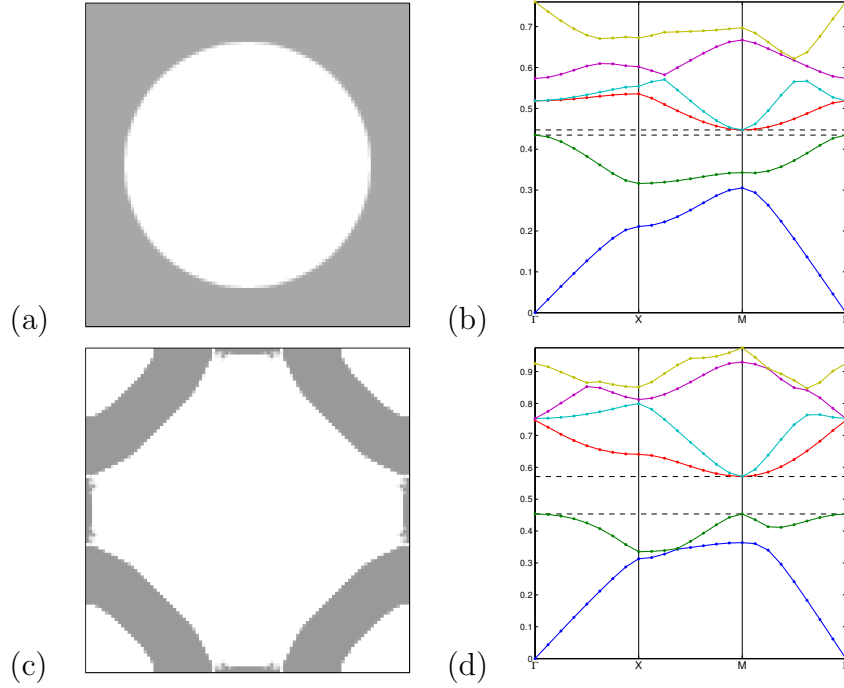


Figure 9.10: Results of Example 9.7. (a) Initial permittivity. (b) Initial band structure. (c) Final permittivity after 87 iterations. (d) Final band structure.

discrete permittivity $\varepsilon_{r,h}^{(58)}$ and TE band structure are depicted in Figure 9.9(c) and (d). The maximal and minimal function values of the final discrete permittivity $\varepsilon_{r,h}^{(58)}$ were equal to 9 and 1, respectively. We also found that the range of $\varepsilon_{r,h}^{(58)}$ was a subset of $\Xi(10^{-5})$, and that $\varepsilon_{r,h}^{(58)}$ took values in $\Xi(10^{-6})$ at 99.9% of all quadrature points. Compared to the Examples 1 and 2, this result is significantly better. Moreover, we found that $|\tilde{w}_{1,h} - w_{1,h}^{(58)}|/\tilde{w}_1 < 10^{-9}$, where \tilde{w}_1 denotes the approximate gap width for the two-valued relative electric permittivity $\tilde{\varepsilon}_{r,h} = \tilde{H}_h(\varepsilon_{r,h}^{(58)})$.

Example 9.7. As can be seen in Figure 9.9, the initial TE band structure of the photonic crystal studied in Example 9.6 also exhibits a band gap between the second and third band with an approximate gap width of $w_{2,h} = 0.0243 \cdot 2\pi c_0/a$. By initializing Algorithm 7.1 exactly as in Example 9.6, we obtain a reduced initial gap width of $w_{2,h}^{(0)} = 0.0126 \cdot 2\pi c_0/a$ (see Figure 9.10(a) and (b)). After only 87 iterations the algorithm converged to a final discrete permittivity $\varepsilon_{r,h}^{(87)}$, which yielded a gap width of $w_{2,h}^{(87)} = 0.1174 \cdot 2\pi c_0/a$ (see Figure 9.10(c) and (d)). We found the maximal and minimal function values of $\varepsilon_{r,h}^{(87)}$ to be 9 and 1, respectively. Moreover, $\varepsilon_{r,h}^{(87)}$ took function values in $\Xi(10^{-4})$ at 99.9% of

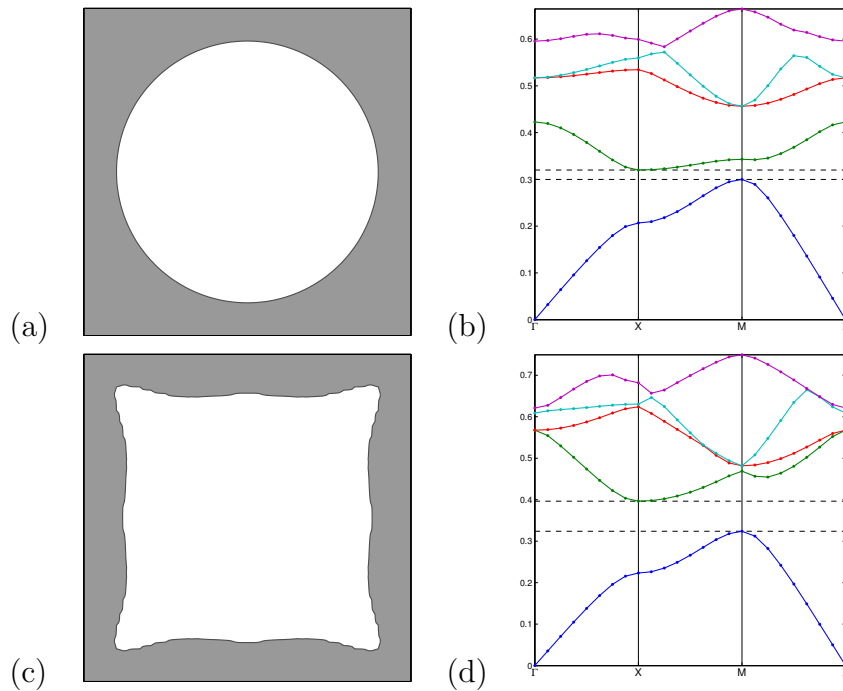


Figure 9.11: Results of Example 9.8. (a) Initial permittivity. (b) Initial band structure. (c) Final permittivity after 4 iterations. (d) Final band structure.

all quadrature points. In view of these results, one is surprised to find that the two-valued permittivity $\tilde{\varepsilon}_{r,h} = \tilde{H}_h(\varepsilon_{r,h}^{(87)})$ yields a smaller gap width of $\tilde{w}_2 = 0.1170 \cdot 2\pi c_0/a$. We attribute this rather large deviation to fact that the final photonic crystal structure in some parts shows cusp-like medium interfaces, which cannot be resolved properly by a uniform finite element mesh.

9.2.2 Results of the Level Set Method

Here we also present some results obtained by Algorithm 8.1. The algorithm parameters were chosen as $t_{\min} = 0.01$, $t_{\max} = 0.1$, $\Delta_{\min} = 10^{-8}$. Descent directions were chosen by a choice strategy similar to that described in Section 7.6.

Example 9.8. In our first example we start off with an initial photonic crystal which is similar to that in Example 9.6. The medium structure of this crystal is described by the initial level set function

$$\varphi^{(0)}(\mathbf{x}) := |\mathbf{x}| - 0.4a,$$

where a denotes the lattice constant. Our aim is to maximize the band gap between the first and the second TE band. We computed an approximate initial gap

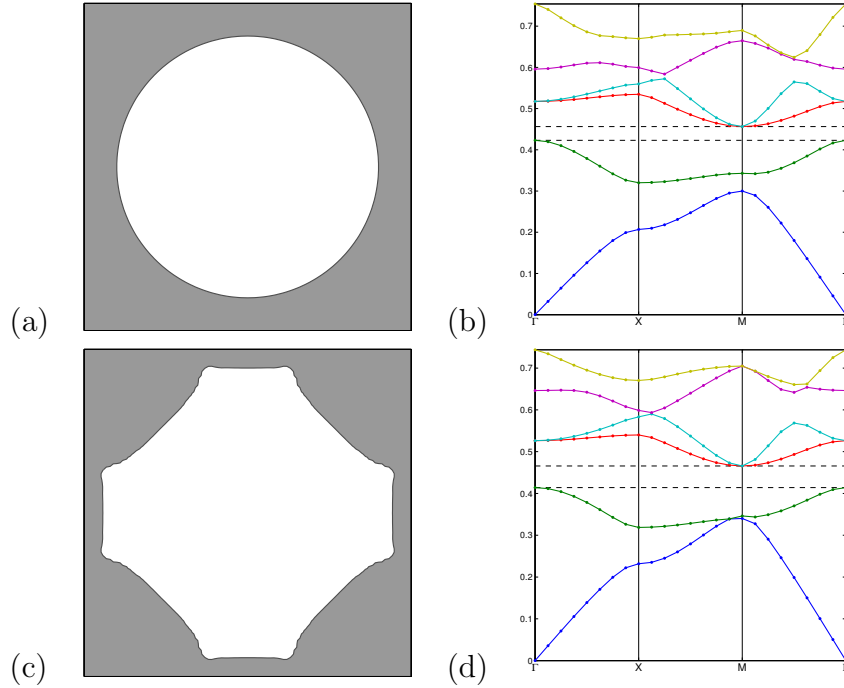


Figure 9.12: Results of Example 9.9. (a) Initial permittivity. (b) Initial band structure. (c) Final permittivity after 4 iterations. (d) Final band structure.

width of $w_{1,h}^{(0)} = 0.0201 \cdot 2\pi c_0/a$ (see Figure 9.11(a) and (b)). The level set method stopped after only 4 iterations yielding a band gap on $w_{1,h}^{(4)} = 0.0726 \cdot 2\pi c_0/a$ (see Figure 9.11(c) and (d)). When we compare this result to the result found by the generalized gradient method in Example 9.6, where a gapwidth of $0.0771 \cdot 2\pi c_0/a$ was achieved, we are lead to conclude that the level set method stopped prematurely. By comparing the final medium structures, we notice that the optimal structure exhibits cusp-shaped material interfaces. We believe that these interfaces present a problem to the level set method.

Example 9.9. In this last example of a band gap maximization problem for two dimensional crystals, we start with the same initial photonic crystal as in the previous example, but aim at maximizing the band gap between the second and the third TE band. The initial gap width is given by $w_{2,h}^{(0)} = 0.0334 \cdot 2\pi c_0/a$ (see Figure 9.12(a) and (b)). As in the previous example the level set method stopped after 4 iterations failing to converge to an optimal solution. The gap is only widened to a final gap width of $w_{2,h}^{(4)} = 0.0517 \cdot 2\pi c_0/a$ (see Figure 9.12(c) and (d)). Compared to the result in Example 9.7, where a final gap width of $0.1174 \cdot 2\pi c_0/a$ could be achieved, the result of the level set method is far from

optimal. When comparing the final medium structures, one notices that they obviously have different topologies. Apparently, the level set method was not able to change the medium structure's topology and thus converged to a less optimal result.

9.3 Maximizing Band Gaps of a 3D Photonic Crystal

In this section we present numerical solutions of photonic band gap maximization problems for three-dimensional photonic crystals. The results we present here are novel and, to the best of our knowledge, no similar results have been published yet. The results present in this section were obtained by the generalized gradient method as given by Algorithm 7.1.

The band gap maximization problems in the examples below are given by the minimization problem

$$\underset{\rho \in \mathcal{C}}{\text{minimize}} \quad J_2(\rho),$$

i.e., we aim at maximizing the band gap between the second and the third photonic bands. In both examples the initial medium structure of the photonic crystal is given by a scaffold-like structure. The crystallographic space group of this structure is isomorphic to the semidirect product $\mathbb{Z}^3 \rtimes O_3(\mathbb{Z})$. The crystal's Bravais lattice is the simple cubic lattice $a\mathbb{Z}^3$, where a denotes the lattice constant. The width of the scaffold beams is given by $0.25a$.

The relative electric permittivity of the beams was given by $\varepsilon_{r,\max} = 13$, that of the ambient medium by $\varepsilon_{r,\min} = 1$. The admissible set was hence given by

$$\mathcal{C} := \left\{ \rho \in \mathcal{E} \mid \rho = \frac{1}{13} + \frac{12}{13} \chi_S, \quad S \in \mathcal{M}(\Omega) \right\}.$$

The problems were discretized by an $\mathbf{H}(\text{curl})$ -conforming finite element ansatz on a hexahedral mesh. We used a quadrature rule with 9 quadrature points on each element. The finite element discretization was implemented in C++ using the parallel finite element library M++ developed by Wieners (cf. [76], [77]). The approximate band structures were computed using a parallel, iterative eigenvalue solver developed by A. Buloyatov (cf. [17]). We used the C++ library `lp_solve` to solve the linear programs arising from the choice strategy for the descent directions (see Section 7.6).

Example 9.10. In this first example the primitive domain Ω was discretized by a $32 \times 32 \times 32$ mesh. We computed an initial gap width of $w_{2,h}^{(0)} = 0.0953 \cdot 2\pi c_0/a$ (see Figure 9.13(a) and (b)). The generalized gradient method stopped after 26

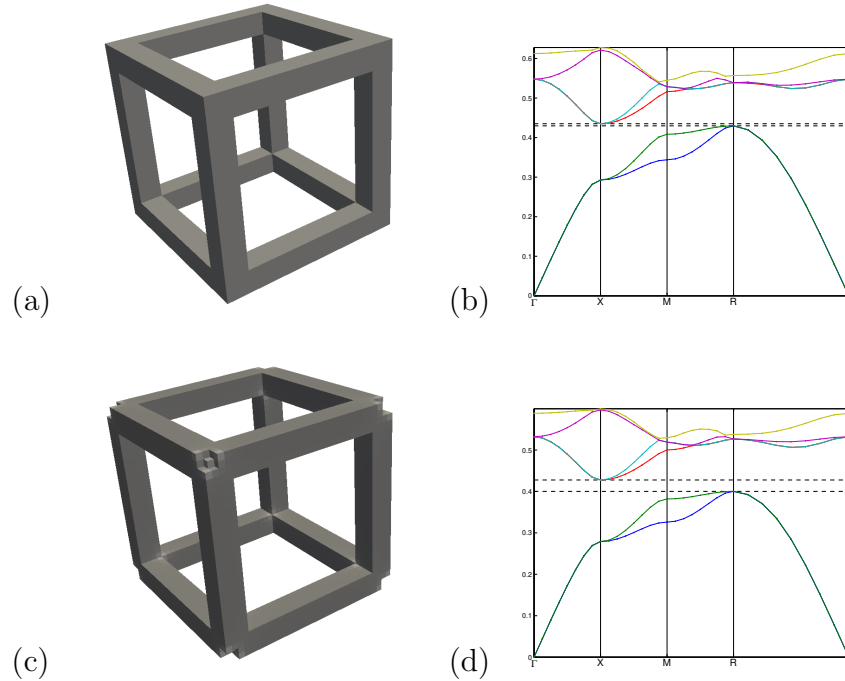


Figure 9.13: Results of Example 9.10. (a) Initial permittivity. (b) Initial band structure. (c) Final permittivity after 26 iterations. (d) Final band structure.

iterations. The final gap width was given by $w_{2,h}^{(26)} = 0.1721 \cdot 2\pi c_0/a$. The minimal and maximal function values of the final discretized relative electric permittivity $\varepsilon_{r,h}^{(26)}$ were given by 1.0137 and 12.8821. However, $\varepsilon_{r,h}^{(26)}$ was not close to a two-valued, admissible permittivity function.

Example 9.11. In this second example the primitive domain Ω was discretized by a $64 \times 64 \times 64$ mesh. We computed an approximate initial gap width of $w_{2,h}^{(0)} = 0.0907 \cdot 2\pi c_0/a$ (see Figure 9.14(a) and (b)). After 32 iterations, the generalized gradient method had widened the band gap to a final width of $w_{2,h}^{(32)} = 0.1674 \cdot 2\pi c_0/a$. The minimal and maximal function values of the final discretized relative electric permittivity $\varepsilon_{r,h}^{(32)}$ were given by 1.0137 and 12.8794. As in the previous example, however, $\varepsilon_{r,h}^{(32)}$ was not close to an admissible permittivity function.

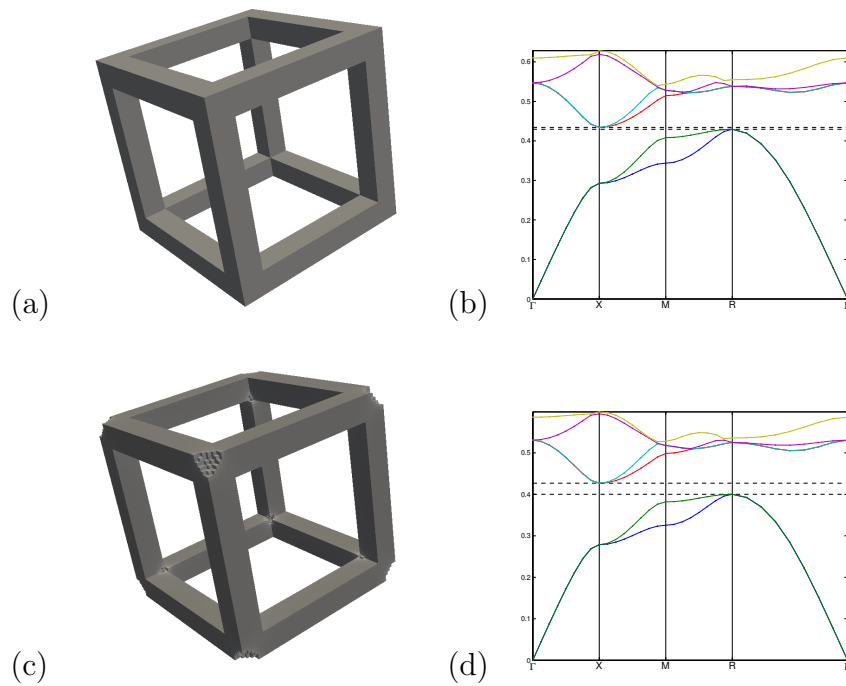


Figure 9.14: Results of Example 9.11. (a) Initial permittivity. (b) Initial band structure. (c) Final permittivity after 32 iterations. (d) Final band structure.

Chapter 10

Summary, Conclusions and Outlook

10.1 Summary

In this dissertation we studied mathematical aspects related to photonic band structure optimization problems and developed optimization algorithms for these problems. The performance of the optimization algorithms was demonstrated by several numerical experiments, included a band gap maximization problem for a three-dimensional photonic crystal.

Based on Maxwell's equations and the Bloch ansatz we derived a family of constrained eigenvalue problems, which determine the existence of so-called Bloch modes in a three-dimensional photonic crystal. Each eigenvalue problem in this family is identified by the so-called quasimomentum vector of the corresponding Bloch modes. The medium structure of the photonic crystal is represented by a coefficient in the eigenvalue equations, which is typically discontinuous. For this reason, we established a weak formulation of the eigenvalue problems. This weak formulation was posed for functions, which belong to certain Sobolev spaces, and which satisfy so-called periodic boundary conditions. The corresponding spectral theory showed that for every photonic crystal and for every quasimomentum vector there exists an increasing, unbounded sequence of non-negative, real eigenvalues. The sequence of corresponding eigenfunctions forms a complete system of the underlying Sobolev space, which is orthogonal in the L^2 -sense.

Given a specific quasimomentum vector, each eigensolution of the corresponding eigenvalue problem determines a Bloch mode, which is able to propagate in the photonic crystal. The eigenvalue determines the phase frequency of such a Bloch mode. We showed that the Bloch mode frequencies depend continuously on the quasimomentum vector. The graphs of the Bloch mode frequencies with respect to

the quasimomentum vector are called the photonic bands of the photonic crystals. The union of all photonic bands is referred to as the photonic band structure of the photonic crystal. We discussed how symmetries in the crystal's medium structure relate to symmetries in its band structure. For certain medium structures, the band structure can exhibit so-called band gaps. In essence, a band gap is a range of frequencies with the property that there exists no Bloch mode, whose phase frequency is in that range.

For certain applications it is desirable that the width of a band gap is large. From this the question arises, how the medium structure of a photonic crystal should be designed, such that its band structure exhibits a band gap with maximal width. The emerging optimization problem is called a photonic band gap maximization problem (PBGMP). More generally, one can formulate so-called photonic band structure optimization problems (PBSOPs), which consist in finding optimal medium structures of a photonic crystal in view of certain features of its band structure. We devised a formal setting for a certain class of PBSOPs. This setting allowed us to treat PBSOPs as minimization problems over an admissible set of two-valued functions. Each function of this admissible set represents a specific medium structure of a photonic crystal. The PBSOP is characterized by a goal functional, which is to be minimized. We showed that the goal functionals are locally Lipschitz continuous. The existence of optimal solutions could only be proved for PBSOPs involving so-called TM band structures of two-dimensional photonic crystals.

In order to devise optimization algorithms for PBSOPs we reviewed the concept of generalized differentials, which extend the concept of differentiability to functionals, which are only locally Lipschitz continuous. By employing known results from nonsmooth analysis, we were able to establish inclusion relations for the generalized differentials of those goal functionals, which define PBGMPs. We were also able to identify supersets of the corresponding generalized gradients.

After commenting on suitable discretization schemes for PBSOPs we developed an iterative optimization algorithm based on the generalized gradients of the discretized goal functionals. The algorithm was designed such that certain optimization constraints, which are inherent to all PBSOPs, are satisfied in each iteration. This was achieved by solving a modified minimization problem on an unconstrained, open admissible set. The performance of the algorithm was demonstrated by several numerical experiments on PBGMPs for two- and three-dimensional photonic crystals. We described a second optimization method, which is based on a level set method. In contrast to the generalized gradient algorithm, the level set algorithm is guaranteed to converge to an admissible, optimal medium structure. We also conducted numerical experiments with the level set algorithm, though only on PBGMPs for two-dimensional photonic crystals.

10.2 Comparison of the Optimization Algorithms

In Chapter 7 and Chapter 8 we developed two optimization algorithms for photonic band structure optimization problems (PBSOPs). One algorithm was based on a generalized gradient method, the other algorithm was based on a level set method. The performance of both algorithms was demonstrated by several numerical experiments on photonic band gap maximization problems (PBGMPs) in Chapter 9. Both algorithms were capable to find optimal medium structures, which maximized certain band gaps in the TM and TE band structures of two-dimensional photonic crystals. The generalized gradient method could also maximize a band gap in the band structure of a three-dimensional photonic crystal. In our experiments we found that the level set method has some significant advantages over the generalized gradient method. One advantage is, that the level set method is by construction guaranteed to converge to a two-valued, optimal solution. In general, the generalized gradient method does not converge to a two-valued solution. Nevertheless, we observed in many experiments that this method converged to an almost two-valued solution. The second clear advantage of the level set method is, that it usually needs less iterations than the generalized gradient method in order to converge to an optimal solution. The computation time needed for a single iteration is about the same for both algorithms. Although level set methods generally cannot be expected to change to topology of a medium structure, we observed that topology changes did occur in some examples. We noticed, however, that the level set method was not able to converge to medium structures with cusp-shaped material interfaces. In that respect, the generalized gradient method performed slightly better.

10.3 Open Problems

In this work we already mentioned some problems, which are still open to the best of our knowledge. Let us briefly review these problems. In Section 5.3, we showed that photonic band structure optimization problems (PBSOPs) involving the TM band structure of a two-dimensional photonic crystal admit an optimal solution within an extended admissible set of functions, which are essentially bounded from above and below. Our numerical experiments suggest that the optimal solutions of photonic band gap maximization problems (PBGMPs) are indeed two-valued functions. However, we could not prove that this is generally the case. We mention here the works of Cox and MacLaughlin [31] and [32], which might be helpful to establish a proof, that optimal solutions of TM PBGMPs are always attained for two-valued relative electric permittivity functions.

For the TE setting, as well as in for three-dimensional setting the existence of optimal solutions could not be proved, either. As we discussed in Section 5.4, one would have to consider the optimization problems on admissible sets, which also contain anisotropic density functions, in order to prove the existence of optimal solutions by standard homogenization methods. The concepts of H-convergence and G-compactness would be central to such a proof. However, we could find conclusive evidence in the literature that these concepts can be transferred to PBSOPs. The question as to whether or not there exist isotropic or two-valued optimal solutions also remained unanswered so far.

Furthermore, we could not verify that the concept of generalized gradient methods can be extended to infinite-dimensional Hilbert spaces, or even to certain Banach spaces such as L^∞ -spaces. The numerical algorithms developed in Chapter 7 mainly rely on the fact that the optimization problems are, in essence, posed on \mathbb{R}^N for some $N \in \mathbb{N}$. This is due to the discretization of PBSOPs. Notice however, that the discretization acts as a regularization operation. Therefore, the generalized gradient method in Chapter 7 only solves regularized versions of PBSOPs. We do not know to what extent the solution of a regularized PBSOP may differ from the solution of an unregularized PBSOP.

10.4 Final Remarks and Outlook

One of the aims of this dissertation was to solve photonic band gap maximization problems for three-dimensional photonic crystal. As was shown by the numerical experiments presented in Section 9.3, the generalized gradient method is able to widen the band gap in a three-dimensional photonic crystal. The main problem with this method is, however, that it does not converge to two-valued density functions in general. In contrast to this, the level set method is guaranteed to converge to two-valued solutions. Therefore, we plan to implement and test this method also for three-dimensional photonic crystals.

Appendix A

A FEM Toolbox for MATLAB

In the course of our research project we developed a finite element toolbox for MATLAB. This toolbox, which we shall refer to as the *FEM Toolbox* hereafter, was used to discretize all photonic band gap maximization problems (PBGMPs) involving TM and TE band structures of two-dimensional photonic crystals. In this appendix chapter we give a brief overview over this toolbox. In Section A.1 we comment on general aspects related to the FEM Toolbox, such as design principles and intended use. In Section A.2 we briefly describe some fundamental data structures and algorithms which are provided by the FEM Toolbox. In Section A.3 we present some features of the FEM Toolbox, which were designed specifically for the numerical solution of PBGMPs. The chapter is completed with an implementation example in Section A.4.

A.1 General Remarks

The FEM Toolbox was developed with the intent to provide an easy-to-use and easy-to-adapt MATLAB function library, which can be used to implement finite element methods in one and two space dimensions. The FEM Toolbox can be used either as a tool for numerical computations or as a basic framework for research and education related to finite element methods. Since MATLAB provides a wide range of standard numerical algorithms and advanced visualization capabilities, we chose to implement the FEM Toolbox entirely in the MATLAB programming language. Wherever possible, mathematical entities are represented by MATLAB arrays or by MATLAB structure arrays. Object-oriented programming concepts, such as MATLAB classes, are not used due to their slower performance.

The development of the FEM Toolbox was mostly motivated by our research on photonic crystals. In order to compute band structures of one and two-dimensional photonic crystals, we needed a finite element software, which supported H^1 -con-

forming one-dimensional and preferably quadrilateral elements, and which could handle periodic boundary conditions. We also required that the software could easily be adapted to the algorithmic requirements of band structure optimization algorithms.

In its current version the FEM Toolbox is capable to discretize boundary value problems, which are of the form

$$\begin{cases} -\nabla \cdot (\boldsymbol{\kappa} \nabla u) + \boldsymbol{\gamma} \cdot \nabla u + \mu u = f & \text{in } \Omega, \\ u - \sigma(u \circ \boldsymbol{\iota}) = r & \text{on } \Gamma_0, \\ -\boldsymbol{\nu} \cdot (\boldsymbol{\kappa} \nabla u) + \eta u = g & \text{on } \Gamma_1 \end{cases} \quad (\text{A.1})$$

for the unknown function $u \in H^1(\Omega)$. Here, Ω denotes a computational domain in \mathbb{R}^n , where $n \in \{1, 2\}$, which is bounded by a Lipschitz smooth curve $\Gamma = \partial\Omega$. The boundary Γ is partitioned into two part Γ_0 and Γ_1 , where either boundary part can be empty. In the above system of equations, $\boldsymbol{\nu}$ denotes the outer unit normal field on Γ_1 . The mapping $\boldsymbol{\iota} : \Gamma_0 \rightarrow \Gamma_0$ identifies certain points on Γ_0 .

In the partial differential equation $\boldsymbol{\kappa} : \Omega \rightarrow \mathbb{R}^{n \times n}$, $\boldsymbol{\gamma} : \Omega \rightarrow \mathbb{R}^n$ and $\mu : \Omega \rightarrow \mathbb{R}$ are given coefficients. The right-hand side of the equation is given by a function $f : \Omega \rightarrow \mathbb{R}$. The equation posed on Γ_0 is referred to as the *essential boundary condition* of the problem. Depending on the choice of $\sigma : \Gamma_0 \rightarrow \mathbb{R}$ and $r : \Gamma_0 \rightarrow \mathbb{R}$, the essential boundary condition is a Dirichlet-type condition (for $\sigma = 0$), periodic boundary condition (for $\sigma = 1$ and $r = 0$) or a quasi-periodic boundary condition (for $r = 0$), provided that $\boldsymbol{\iota}$ is chosen appropriately. The equation posed on Γ_1 is referred to as the *natural boundary condition* of the problem. Depending on the choice of $\eta : \Gamma_1 \rightarrow \mathbb{R}$ and $g : \Gamma_1 \rightarrow \mathbb{R}$, this condition is a Neumann-type condition (for $\eta = 0$) or a Robin-type condition (for $\eta \neq 0$). In the next section we briefly describe, how the different entities are represented by the FEM Toolbox, and how the discretization of a boundary value problem of the above form (A.1) is realized. We remark that the FEM Toolbox was solely used to produce the numerical results presented in Chapter 9 in this work, as well as the numerical results presented in [36].

It should be noted that MATLAB's *Partial Differential Equations Toolbox* also provides data structures and algorithms for the implementation of finite element methods. In particular, this toolbox provides sophisticated concepts to describe the geometry of computational domains, as well as a mesh generator. Moreover, it features a graphical user interface which can be used to define and solve various types of elliptic boundary value problems. The reason why we did not use this toolbox was that it only provides data structures and algorithms for H^1 -conforming, linear finite elements on triangular meshes. Furthermore, the Partial Differential Equations Toolbox is not able to handle periodic boundary conditions.

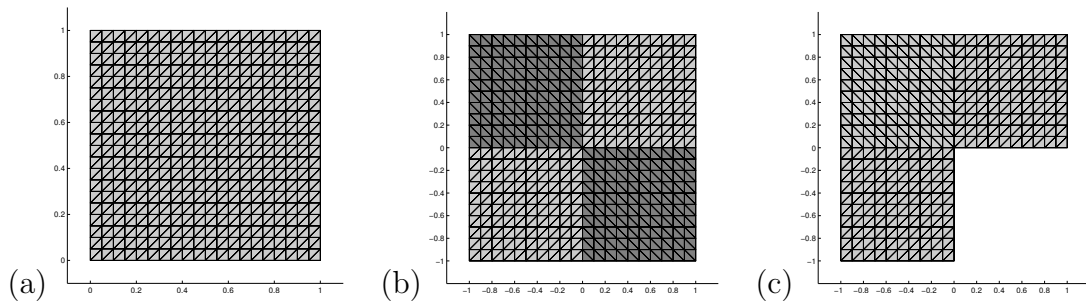


Figure A.1: Finite element meshes returned by the functions `femsquarem` (a), `femcheckerm` (b), and `femlshapem` (c). Different shades of gray indicate different subdomains.

A.2 Data Structures and Algorithms

In this section we describe some important data structures and algorithms, which are provided by the FEM Toolbox. The section does not give a complete account of the toolbox’s functionalities and only introduces the most basic concepts.

The fundamental data structure of the FEM Toolbox is a MATLAB structure array, which describes the mesh of a finite element discretization. This *mesh structure* is capable to represent one-dimensional, triangular, quadrilateral or hybrid, triangular-quadrilateral meshes. The mesh structure stores information about the mesh’s vertices, edges and faces. In addition to that, the *mesh structure* also stores information about boundary vertices and boundary edges. In particular, the mesh structure stores information about boundary identifications, which were represented by the mapping ι in the previous section. Furthermore, the mesh structure stores a *subdomain index* for every mesh element, and a *boundary part index* for every boundary edge or vertex. With this, a computational domain can be subdivided into several *subdomains*, and its boundary can be subdivided into several *boundary parts*.

The FEM Toolbox provides a collection of MATLAB functions that create mesh structures representing specific finite element meshes. Here, we only mention the functions `feminterm`, `femsquarem`, `femcheckerm` and `femlshapem`, which create mesh structures that represent finite element meshes on an interval, on a square, on a 2×2 checkerboard, and on an L-shaped domain, respectively (see Figure A.1). The different “fields” of the checkerboard are distinguished by the subdomain indices mentioned above. As a convention, all MATLAB functions of the FEM Toolbox begin with the letters “fem”. The FEM Toolbox provides a number of functions, which can be used to retrieve informations about a given mesh. Given a mesh structure `msh`, the function call `femnelems(msh)`, for example, re-

turns the number of elements in mesh. Other functions, such as `femnsubdoms`, `femnbparts`, or `femhmax` return the number of different subdomains, the number of different boundary parts, or the maximal element diameter in a mesh. One-dimensional and triangular meshes can also be refined, either globally or locally, by the function `femrefine`.

Finite element spaces are represented by so-called *space structures*. A space structure is a MATLAB structure array, which stores the degrees of freedoms on every *node* in a finite element mesh. Here, the term “node” refers to a vertex, an edge, or a face in a finite element mesh. The space structure also stores informations about the shape functions of the finite element space. In the current version of the FEM Toolbox, space structures are used to represent higher-order H^1 - and L^2 -conforming finite element spaces. We implemented Lobatto-type finite element shape functions, as defined in [73].

Finite element spaces can be created on given meshes by specific MATLAB functions, which are provided by the FEM Toolbox. Given a mesh structure `msh`, the function call `femh1(msh)`, for example creates a space structure, which represents a linear, H^1 -conforming finite element space on the mesh, which is represented by `msh`. Other MATLAB functions provided by the FEM Toolbox retrieve informations about finite element spaces. Given a space structure `spc`, the function call `femndofs(spc)`, for example, returns the number of degrees of freedom in the finite element space represented by `spc`. Finite element spaces can be h-refined or p-refined by the functions `femhrefine`, `femprefine`. Both function can be used to perform either global or local refinements. The function `feminterp` interpolates a MATLAB function on a finite element space and returns the corresponding *coefficient vector*.

The partial differential equation

$$-\nabla \cdot (\boldsymbol{\kappa} \nabla u) + \boldsymbol{\gamma} \cdot \nabla u + \mu u = f \quad \text{on } \Omega$$

of a boundary value problem of the form (A.1) is represented by a MATLAB structure array, which is called the *PDE structure*. A PDE structure is typically generated by the FEM Toolbox function `fempde`, whose arguments define the coefficients $\boldsymbol{\kappa}$, $\boldsymbol{\gamma}$, μ , as well as the right-hand side f of an elliptic partial differential equation. The coefficients and the right-hand side can each be specified subdomain-wise, element-wise or quadrature-point-wise.

Boundary conditions of the form

$$\begin{aligned} u - \sigma(u \circ \boldsymbol{\nu}) &= r && \text{on } \Gamma_0, \\ -\boldsymbol{\nu} \cdot (\boldsymbol{\kappa} \nabla u) + \eta u &= g && \text{on } \Gamma_1 \end{aligned}$$

are represented by so-called *boundary condition structures*. These MATLAB structure arrays, in essence, store informations about the type of boundary condition,

as well as the boundary data. The boundary data r , σ , r , η , and g can be specified boundary-part-wise. Several standard boundary conditions can be defined for a given finite element space by MATLAB functions, which are also provided by the FEM Toolbox. Given a mesh structure `msh` and a space structure `spc`, such that `spc` represents a finite element space on the mesh which is represented by `msh`, the function call `femperiodic(msh,spc)` generates a boundary condition structure representing periodic boundary conditions. Other functions, which generate boundary condition structures are `femdirichlet` or `femneumann`.

After defining a mesh, a finite element space on the mesh, and a partial differential equation, one can assemble the corresponding finite element matrices using the function `femassemble`. This function returns sparse matrices `K`, `C`, and `M`, as well as a vector `f`, which are the *stiffness matrix*, the *damping matrix*, the *mass matrix*, and the *load vector*.

In order to incorporate the boundary conditions, one also needs to assemble the matrices `E` and `H`, as well as the vectors `r` and `g` using the function `femassembcs`. The matrix `E` and the vector `r` represent essential boundary conditions. Recall that essential boundary conditions reduce the number of degrees of freedom in a finite element space and thus give rise to a constrained finite element space. The matrix `E` realizes the identical embedding from the constrained finite element space into the unconstrained finite element space. The vector `r` represents the boundary data. Given a coefficient vector `u` in the constrained finite element space, one obtains the corresponding coefficient vector `v` in the unconstrained space by

$$\mathbf{v} = \mathbf{E} * \mathbf{u} + \mathbf{r};$$

The matrix `H` and the vector `g` represent natural boundary conditions. Natural boundary conditions modify the left-hand side, as well as the right-hand side in the linear system associated with the boundary value problem. In order to obtain this linear system, one needs to compute *system matrix* `A` and the *right-hand side vector* `b` according to

$$\begin{aligned} \mathbf{A} &= \mathbf{E}' * (\mathbf{A} + \mathbf{H}) * \mathbf{E}; \\ \mathbf{b} &= \mathbf{E}' * (\mathbf{f} + \mathbf{g} - (\mathbf{A} + \mathbf{H}) * \mathbf{r}); \end{aligned}$$

The matrices `A` and `E`, as well as the vectors `b` and `r` can also be computed directly by the FEM Toolbox function `femassembvp`. The coefficient vector `u` of an approximate solution of the boundary value problem can then be computed using MATLAB's standard solver routines such as the MATLAB backslash `\`, `pcg`, or `bicg`.

Eigenvalue problems can be discretized in a similar way. After assembling the matrices `K`, `M`, and `E` mentioned above one obtains the left-hand side matrix `A` and the right-hand side matrix `B` of the associated generalized matrix eigenvalue problem by

$$\begin{aligned} \mathbf{A} &= \mathbf{E}' * \mathbf{K} * \mathbf{E}; \\ \mathbf{B} &= \mathbf{E}' * \mathbf{M} * \mathbf{E}; \end{aligned}$$

Alternatively, one can assemble the matrices \mathbf{A} , \mathbf{B} , and \mathbf{E} directly by using the FEM toolbox function `femassemvp`. In order to obtain approximate eigenvalues as well as the coefficient vectors of approximate eigenfunctions, one can solve the generalized matrix eigenvalue problem by using MATLAB's `eigs` function.

Most of the FEM Toolbox's functionalities are illustrated by a number of demo programs. These demo programs are MATLAB scripts which can be started through the commands `femdemo0`, `femdemo1`, ..., `femdemo12` (in the current version). In particular, the demo program `femdemo0` provides a "crash course" for users, who simply want to use the FEM Toolbox in order to solve a standard elliptic boundary value problem numerically. All functions in the FEM Toolbox are documented according to MATLAB conventions. A user manual is not available yet.

A.3 Some Customized Features

In this section we present some features of the FEM Toolbox, which were developed specifically for the optimization algorithms introduced in Chapter 7 and Chapter 8.

In the generalized gradient method, as well as in the level set method, one often needs to evaluate the approximate eigenfunctions at the quadrature points in the mesh. Clearly, point evaluations are linear operations and can therefore be represented by matrices. We hence added a function named `femqdelta` to the FEM Toolbox. This function assembles a sparse matrix \mathbf{D} , which evaluates finite element functions at all quadrature points in the corresponding finite element mesh. Given a coefficient vector \mathbf{u} , one obtains by

$$\mathbf{y} = \mathbf{D} * \mathbf{u};$$

a vector \mathbf{y} , which contains the function values of the finite element function represented by \mathbf{u} at the quadrature points in the mesh. Thus, point evaluations at quadrature points can be performed by a single matrix-vector product. As long as the finite element mesh remains unchanged, which was the case for the above mentioned optimization algorithms, the matrix \mathbf{D} does not need to be reassembled.

We found that the matrix \mathbf{D} could also be used to speed up reassemblies of finite element matrices. Recall that the discretized density function ρ changed in each iteration of the optimization algorithms for photonic band structure optimization problems. As a consequence, certain finite element matrices needed to be reassembled in order to compute the band structure for the new medium structure. In the TM setting, the mass matrix \mathbf{M} depends on the coefficient ρ and needs to

be reassembled after every iteration. In order to speed up this reassembly, we devised the function `fempreassem`, which in particular returns a sparse diagonal matrix W , such that the mass matrix for the coefficient equal to 1 can be computed by

$$M = D' * W * D;$$

The matrix product on the right-hand side realizes the assembly process for the mass matrix M . The elements on the diagonal of W are precisely the weights of the quadrature rule, which is used during the assembly. A fast reassembly of the mass matrix can now be performed as follows. Suppose that `rho` is a vector, which contains the function values of the discretized density function ρ at all quadrature points in the mesh. One can then construct a sparse diagonal matrix `Rho` with the elements of `rho` on the diagonal. By

$$Mrho = D' * W * Rho * D;$$

one obtains the mass matrix `Mrho` for the coefficient ρ . As we expected, evaluating the above matrix product was significantly faster than a standard matrix assembly. The function `fempreassem` returns further matrices, which can be used to assemble the stiffness and the damping matrix also by evaluating specific matrix products.

A.4 Implementation Examples

In this final section, we demonstrate the use of the FEM Toolbox through small implementation examples. In the first example, our aim is to compute approximations of the four smallest eigenvalues $\lambda_1, \dots, \lambda_4$ of the Laplace eigenvalue problem

$$\begin{cases} -\Delta u = \lambda u & \text{in } \Omega, \\ u = 0 & \text{on } \partial\Omega, \end{cases}$$

where Ω denotes the square $(0, \pi)^2$. Clearly, the eigenvalues of the above problem are given by $\lambda^{(m,n)} := m^2 + n^2$, where $m, n \in \mathbb{N}$. Hence, the four smallest eigenvalues are given by $\lambda_1 := \lambda^{(1,1)} = 2$, $\lambda_2 := \lambda^{(2,1)} = 5$, $\lambda_3 := \lambda^{(1,2)} = 5$, and $\lambda_4 := \lambda^{(2,2)} = 8$. Notice that the number 5 is a double eigenvalue.

The MATLAB code listed as Program A.1 demonstrates, how this problem can be discretized with the FEM Toolbox. At first, a mesh structure `msh` is created. The mesh structure represents a 20×20 structured, quadrilateral mesh on Ω . Next, an H^1 -conforming, bilinear finite element space is created on the mesh. The finite element space is represented by the space structure `spc`. The PDE structures `lsh`

Program A.1 Approximation of the four smallest Laplace eigenvalues

```

% Create a 20-by-20 rectangular structured mesh
% on a square with side lengths equal to 2*pi.
n = 20;
msh = femsquarem(n,0,0,pi,'rectangular');

% Create an H1-conforming, linear finite element
% space on the above mesh.
p = 1;
spc = femh1(msh,p);

% Define the left-hand side and the right-hand side
% operators of the standard Laplace eigenvalue equation.
lhs = fempde(msh,spc,1);
rhs = fempde(msh,spc,0,0,1);

% Define homogeneous Dirichlet conditions.
bcs = femdirichlet(msh,spc);

% Assemble the finite element matrices for the generalized
% matrix eigenvalue problem  $A * v = \lambda * M * v$ .
[A,B,E] = femassemvp(msh,spc,bcs,lhs,rhs);

% Compute the four smallest eigenvalues
k = 4;
s = 0;
[V,Lambda] = eigs(A,B,k,s);

% Realize the boundary conditions
U = E * V;

% Plot the eigenfunction
for i = 1:4
    figure(i);
    femplot(msh,spc,U(:,5-i));
end

```

and `rhs` represent the operator $-\Delta$ ($\kappa = 1$, $\gamma = 0$, $\mu = 0$) and the identity operator ($\kappa = 0$, $\gamma = 0$, $\mu = 1$), respectively. The boundary condition structure `bcs` represents the homogeneous Dirichlet boundary conditions imposed on $\partial\Omega$. The

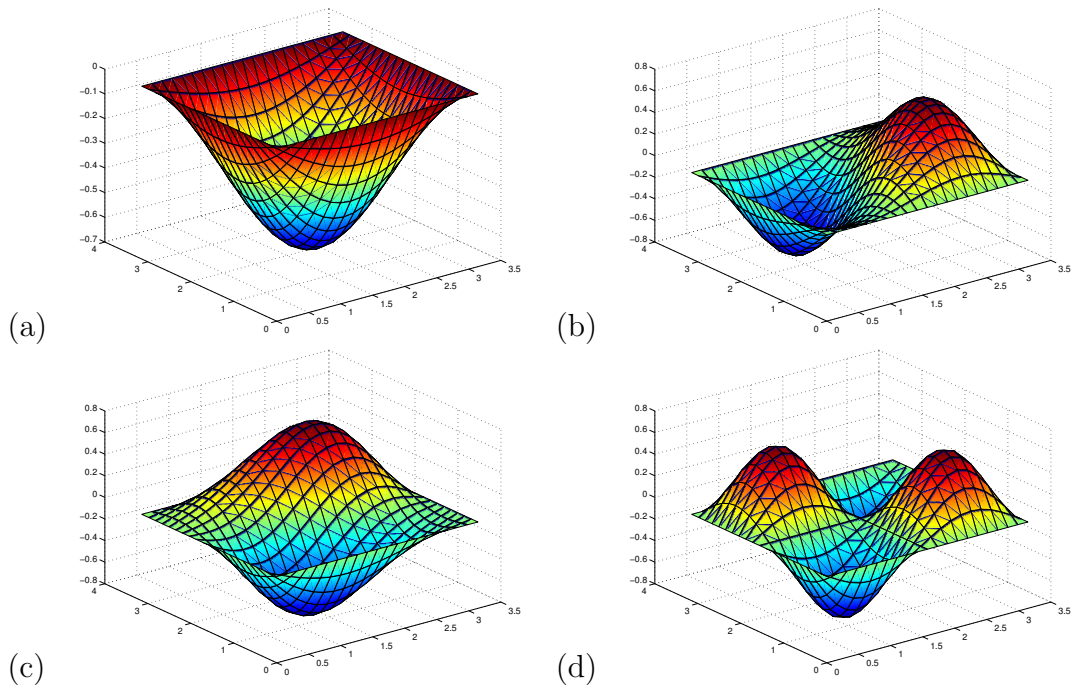


Figure A.2: Eigenfunctions corresponding to the four smallest Laplace eigenvalues. The plots (b) and (c) show eigenfunctions that correspond to the double eigenvalue λ_2 .

finite element matrices are assembled using the function `femassemevp`. Approximations for the four smallest eigenvalues are then computed by MATLAB's `eigs` function. The function `eigs` also returns the coefficient vectors of the corresponding approximate eigenfunctions. However, these coefficient vectors only determine the eigenfunction's interior degrees of freedom. The degrees of freedom at the boundary are realized by virtue of the boundary conditions. In the code, this realization is accomplished by applying the embedding matrix \mathbf{E} to the coefficient vectors. Finally, the eigenfunctions are plotted. The plots are shown in Figure A.2.

In our second example we investigate the convergence properties of the eigenvalue approximations with respect to the mesh width h_{\max} , i.e., with respect to the largest element diameter in the finite element mesh. Consider the MATLAB code listed as Program A.2. At first, a coarse triangular mesh on the computational domain Ω is created. An H^1 -conforming finite element space on the mesh is created next. Notice that this space consists of linear finite elements, only (`p = 1`). In the following loop the mesh, as well as finite element space are repeatedly globally refined. The refinement is accomplished by a newest-node bisection algorithm (see e.g. Chapter 3 in [68]). In Figure A.3(a) we depicted the first three refined

Program A.2 Convergence analysis of the four smallest Laplace eigenvalues

```

% Define a coarse finite element discretization
p = 1;
msh = femsquarem(4,0,0,pi,'criss');
spc = femh1(msh,p);

% Allocate memory
hmax = zeros(1,6);
error = zeros(4,6);

for j = 1:8
    % Refine the finite element space
    [msh,spc] = femhrefine(msh,spc);

    % Mesh width and number of DOFs
    hmax(j) = femhmax(msh);

    % Solve the boundary value problem
    lhs = fempde(msh,spc,1);
    rhs = fempde(msh,spc,0,0,1);
    bcs = femdirichlet(msh,spc);

    % Compute the four smallest eigenvalues
    [A,B] = femassemvp(msh,spc,bcs,lhs,rhs);
    lambda = eigs(A,B,4,0);

    % Determine the errors
    error(:,j) = abs([8; 5; 5; 2] - lambda);
end

```

meshes. After every refinement, the current mesh size h_{\max} is stored in the array `hmax`. Then, approximations of the smallest four Laplace eigenvalues on Ω are computed. The absolute approximation errors are then stored in the columns of the matrix `error`.

First, we ran Program A.2 as listed with linear finite element spaces. Afterwards, we ran the same program with quadratic finite element spaces, i.e., we replaced the statement `p = 1;` by `p = 2;`. The results of both program runs are displayed in Figure A.3(b) and (c). As expected, we observed an experimental order of convergence (EOC) of approximately 2 for linear finite elements, and an EOC of approximately 4 for quadratic finite elements. We remark that the EOCs

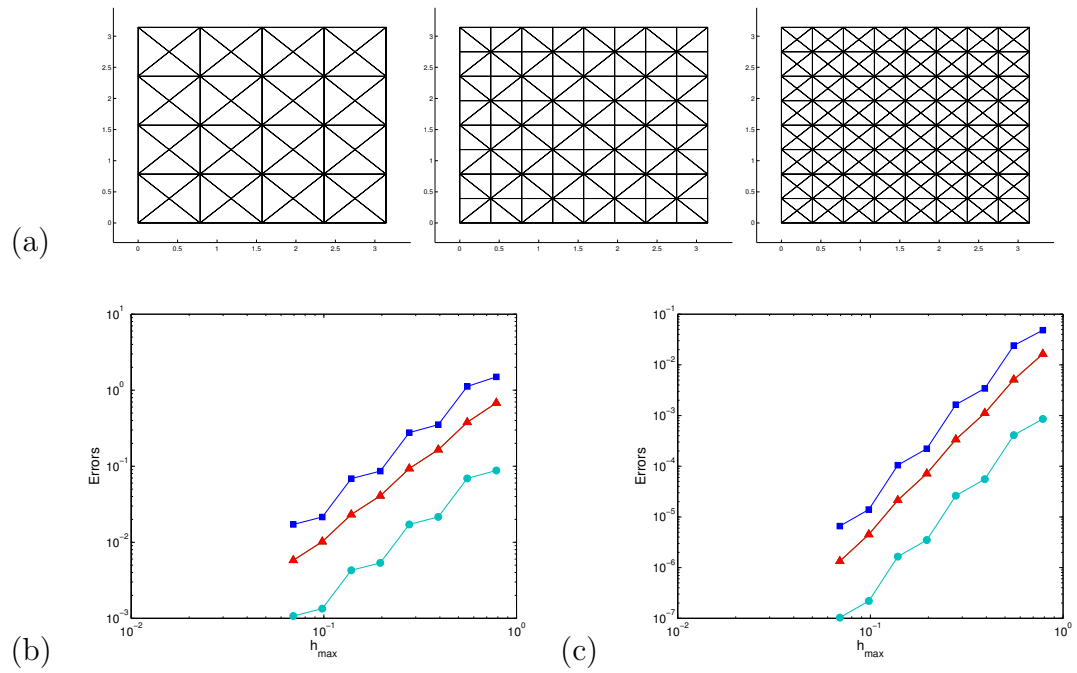


Figure A.3: The first three refined mesh generated in Program A.2 (a) and the absolute approximation errors for linear elements (b) and for quadratic elements (c). Round markers (●) indicate errors for the eigenvalue λ_1 , triangular markers (▲) indicate errors for the eigenvalue $\lambda_2 = \lambda_3$, and square markers (■) indicate errors for the eigenvalue λ_4 .

can be computed with the function `femheoc`, which is also provided by the FEM Toolbox.

Frequently Used Symbols

$A_{\mathbf{k}}(\rho)$	The conjugate-linear operator induced by $a_{\mathbf{k}}(\rho)$, see (4.43) on page 65
$a_{\mathbf{k}}(\rho)$	Sesquilinear form on $\mathbf{W} \times \mathbf{W}$, see (4.33) on page 60
$a_{\mathbf{k}}^{\text{TE}}(\rho)$	Sesquilinear form on $V \times V$, see (4.92) on page 87
$a_{\mathbf{k}}^{\text{TM}}$	Sesquilinear form on $V \times V$, see (4.81) on page 85
\mathbb{B}	The first Brillouin zone of Λ , see (3.34) on page 38
\mathcal{C}	The admissible set, see (5.2) on page 91
$\bar{\mathcal{C}}$	The extended admissible set, see (5.4) on page 95
\mathcal{C}_h	Discretization of \mathcal{C} , see (7.8) on page 132
$\bar{\mathcal{C}}_h$	Discretization of $\bar{\mathcal{C}}$, see (7.9) on page 132
\mathcal{D}	The subset of essentially positive functions in \mathcal{E} , see (4.32) on page 59
\mathcal{D}_h	Discretization of \mathcal{D} , see (7.7) on page 132
\mathcal{E}	The real Banach space $L^\infty(\Omega, \mathbb{R})$, see (4.31) on page 59
\mathcal{E}_h	Discretization of \mathcal{E} , see (7.5) on page 131
$G_{\mathbf{k}}(\rho)$	The Green solution operator for $L_{\mathbf{k}}(\rho)$, see (4.51) on page 66
\mathcal{H}	The real Hilbert space $H_{\text{per}}^1(\Omega, \mathbb{R})$, see (8.1) on page 148
J	The generic goal functional, see (5.3) on page 92
J^{TE}	The generic TE goal functional, see (5.8) on page 104
J^{TM}	The generic TM goal functional, see (5.6) on page 97
$M_{\mathbf{k}}$	The conjugate-linear operator induced by m , see (4.44) on page 65

m	Sesquilinear form on $\mathbf{Z} \times \mathbf{Z}$, see (4.35) on page 60
m^{TE}	Sesquilinear form on $Z \times Z$, see (4.93) on page 87
$m^{\text{TM}}(\rho)$	Sesquilinear form on $Z \times Z$, see (4.82) on page 85
Q	The complex Hilbert space $H_{\text{per}}^1(\Omega)$, see (4.29) on page 59
$\mathbf{u}_j(\rho, \mathbf{k})$	The j -th eigenfunction, see (4.53) on page 68
\mathcal{H}	The real Banach space $W^{1,\infty}(\Omega, \mathbb{R})$, see (8.2) on page 148
V	The complex Hilbert space $H_{\text{per}}^1(\Omega)$, see (4.79) on page 84
$\mathbf{V}_{\mathbf{k}}$	Linear subspace of \mathbf{W} , see (4.36) on page 60
\mathbf{W}	The complex Hilbert space $\mathbf{H}_{\text{per}}(\text{curl}; \Omega)$, see (4.28) on page 59
\mathcal{X}_h	The set of quadrature points, see (7.4) on page 131
Z	The complex Hilbert space $L^2(\Omega)$, see (4.80) on page 84
\mathbf{Z}	The complex Hilbert space $L^2(\Omega)^3$, see (4.30) on page 59
Π_h	The interpolation operator from \mathcal{E} to \mathcal{E}_h , see (7.6) on page 132
Λ	Bravais lattice, see (3.32) on page 38
$\lambda_j(\rho, \mathbf{k})$	The j -th smallest eigenvalue of $A_{\mathbf{k}}(\rho)$, see (4.60) on page 70
$\mu_j(\rho, \mathbf{k})$	The j -th largest eigenvalue of $G_{\mathbf{k}}(\rho)$, see (4.52) on page 68
Ω	Primitive cell of Λ , see (3.36) on page 39
$\omega_j(\rho, \mathbf{k})$	The j -th smallest eigenfrequency, see (4.74) on page 80
$\omega_j^{\text{TE}}(\rho, \mathbf{k})$	The j -th smallest TE eigenfrequency, see (4.96) on page 87
$\omega_j^{\text{TM}}(\rho, \mathbf{k})$	The j -th smallest TM eigenfrequency, see (4.95) on page 87
ρ_{\max}	Essential upper bound for the functions in \mathcal{C} , see (5.1) on page 91
ρ_{\min}	Essential lower bound for the functions in \mathcal{C} , see (5.1) on page 91

About The Author

The author, Markus Thobias Richter, was born on July 8th 1980 in Frankfurt am Main, Germany. After living in Darmstadt–Eberstadt for a few years, his family moved to Alsbach–Hähnlein, a municipality near Darmstadt. There, he attended the local primary school from 1987 onwards. In 1991 he was admitted at the grammar school Altes Kurfürstliches Gymnasium (AKG) in Bensheim, where he received his Abitur in 2000. After fulfilling his Civilian Service as an emergency medical technician, he commenced his undergraduate studies in applied mathematics at the University of Karlsruhe in 2001. There, he moved on to graduate studies in 2003, specializing in the numerical analysis of finite element methods and mathematical models in optical communications technology. After writing his diploma thesis under the supervision of Prof. Dr. Willy Dörfler and Prof. Dr. Christian Wieners on the numerical computation of photonic band structures for dispersive media, he graduated from the University of Karlsruhe in 2006 with the academic degree Diplom–Technomathematiker. In the same year he began his doctoral studies as an associated member of the DFG Research Training Group 1294 at the Department of Mathematics of the University of Karlsruhe. His doctoral advisors were again Prof. Dr. Willy Dörfler and Prof. Dr. Christian Wieners. During his time as a doctoral candidate the author also worked as a research assistant at the Institute of Applied Mathematics and Numerical Analysis at the University of Karlsruhe.

Bibliography

- [1] Robert A. Adams and John J. F. Fournier, *Sobolev Spaces. 2nd Edition*, Pure and Applied Mathematics, no. 140, Academic Press, New York, 2003.
- [2] Grégoire Allaire and François Jouve, *A level-set method for vibration and multiple loads structural optimization*, Comput. Methods Appl. Mech. Eng. **194** (2005), no. 30-33, 3269–3290.
- [3] Grégoire Allaire and Robert V. Kohn, *Optimal design for minimum weight and compliance in plane stress using extremal microstructures*, Eur. J. Mech., A **12** (1993), no. 6, 839–878.
- [4] Hans W. Alt, *Lineare Funktionalanalysis: Eine anwendungsorientierte Einführung. 4. Auflage*, Springer, Berlin, 2002.
- [5] Angelo Alvino, Pierre-Louis Lions, and Guido Trombetti, *On optimization problems with prescribed rearrangements*, Nonlinear Anal., Theory Methods Appl. **13** (1989), no. 2, 185–220.
- [6] Mark A. Armstrong, *Groups and Symmetry*, Springer, Berlin, 1988.
- [7] Douglas N. Arnold, Richard S. Falk, and Ragnar Winther, *Finite element exterior calculus, homological techniques, and applications*, Acta Numerica **15** (2006), 1–155.
- [8] Giles Auchmuty, *Duality for non-convex variational principles*, J. Differ. Equations **50** (1983), 80–145.
- [9] ———, *Dual variational principles for eigenvalue problems*, Nonlinear Functional Analysis and its Applications (Felix E. Browder, ed.), AMS, Providence (RI), 1986, pp. 136–211.
- [10] George Bachman, Lawrence Narici, and Edward Beckenstein, *Fourier and Wavelet Analysis*, Springer, Berlin, 2000.

- [11] Martin P. Bendsøe and Ole Sigmund, *Topology Optimization: Theory Methods and Applications*, Springer, Berlin, 2003.
- [12] Serge Berthier, *Iridescence: The Physical Color of Insects*, Springer, Berlin, 2007.
- [13] Felix Bloch, *Über die Quantenmechanik der Elektronen in Kristallgittern*, Z. Phys. **52** (1928), 555–600.
- [14] Daniele Boffi, Matteo Conforti, and Lucia Gastaldi, *Modified edge finite elements for photonic crystals*, Numer. Math. **105** (2006), no. 2, 249–266.
- [15] Walter Borchartd-Ott, *Crystallography. 2nd edition*, Springer, Berlin, 1995.
- [16] Raoul Bott and Loring W. Tu, *Differential Forms in Algebraic Topology*, Graduate texts in mathematics, no. 82, Springer, Berlin, 1995.
- [17] Alexander Bulovyatov, *A Parallel Multigrid Method for Band Structure Computation of 3D Photonic Crystals with Higher Order Finite Elements*, Ph.D. thesis, Karlsruhe Institute of Technology (KIT), 2010.
- [18] Bureau international des poids et mesures, Sévres, *Le Système international d'unités*, 8 ed., 2006.
- [19] Dan Butnariu and Abraham Mehrez, *Convergence criteria for generalized gradient methods of solving locally lipschitz feasibility problems*, Comput. Optim. Appl. **1** (1992), no. 3, 307–326.
- [20] Michel Cessenat, *Mathematical Methods in Electromagnetism: Linear Theory and Applications*, World Scientific Publishing, Singapore, 1996.
- [21] Doina Cioranescu and Patrizia Donato, *An Introduction to Homogenization*, Oxford University Press, 1999.
- [22] Frank H. Clarke, *Optimization and Nonsmooth Analysis*, Canadian Mathematical Society Series of Monographs and Advanced Texts, Wiley, New York, 1983.
- [23] Frank H. Clarke, Yuri S. Ledyaev, Ron J. Stern, and Peter R. Wolenski, *Nonsmooth Analysis Optimization and Control Theory*, Graduate Texts in Mathematics, no. 178, Springer, Berlin, 1998.
- [24] Philippe Clément, *Approximation by finite element functions using local regularization*, Rev. Franc. Automat. Inform. Rech. Operat. **9** (1975), no. R-2, 77–84.

- [25] Gary C. Cohen, *Higher-Order Numerical Methods for Transient Wave Equations*, Scientific Computation, Springer, Berlin, 2002.
- [26] Carlos Conca, Rajesh Mahadevan, and León Sanz, *An extremal eigenvalue problem for a two-phase conductor in a ball*, Appl. Math. Optim. **60** (2009), no. 2, 173–184.
- [27] Steven Cox and Robert Lipton, *Extremal eigenvalue problems for two-phase conductors*, Arch. Ration. Mech. Anal. **136** (1996), no. 2, 101–117.
- [28] Steven J. Cox, *The generalized gradient at a multiple eigenvalue*, J. Funct. Anal. **133** (1995), no. 1, 30–40.
- [29] Steven J. Cox and David C. Dobson, *Maximizing band gaps in two-dimensional photonic crystals*, SIAM J. Appl. Math. **59** (1999), no. 6, 2108–2120.
- [30] ———, *Band structure optimization of two-dimensional photonic crystals in h-polarization*, J. Comput. Phys. **158** (2000), no. 2, 214–224.
- [31] Steven J. Cox and Joyce R. McLaughlin, *Extremal eigenvalue problems for composite membranes. I*, Appl. Math. Optimization **22** (1990), no. 2, 153–167.
- [32] ———, *Extremal eigenvalue problems for composite membranes. II*, Appl. Math. Optimization **22** (1990), no. 2, 169–187.
- [33] David C. Dobson and Joseph E. Pasciak, *Analysis of an algorithm for computing electromagnetic bloch modes using nedelec spaces*, Comput. Methods Appl. Math. **1** (2001), no. 2, 138–153.
- [34] Georges Duvaut and Jacques-Louis Lions, *Les inéquations en mécanique et en physique*, Dunod, Paris, 1972.
- [35] Ivar Ekeland and Roger Temam, *Analyse Convexe et Problème Variationelles*, Dunod, Paris, 1974.
- [36] Christian Engström and Markus Richter, *On the spectrum of an operator pencil with applications to wave propagation in periodic and frequency dependent materials*, SIAM J. Appl. Math. **70** (2009), no. 1, 231–247.
- [37] Gaston Floquet, *Sur les équations différentielles linéaires à coefficients périodiques*, Ann. de l'Éc. N. **12** (1883), 47–89.

- [38] Vivette Girault and Pierre-Arnaud Raviart, *Finite Element Approximation of the Navier-Stokes Equations*, Lecture Notes in Mathematics, no. 749, Springer, Berlin, 1979.
- [39] Ralf Hiptmair, *Finite elements in computational electromagnetism*, Acta Numerica **11** (2002), 237–339.
- [40] Friedrich Hirzebruch and Winfried Scharlau, *Einführung in die Funktionalanalysis*, Spektrum Akademischer Verlag, Heidelberg, 1971.
- [41] John D. Jackson, *Classical Electrodynamics. 3rd Edition*, John Wiley & Sons, Hoboken (NJ), 1999.
- [42] John D. Joannopoulos, Steven G. Johnson, Joshua N. Winn, and Robert D. Meade, *Photonic Crystals: Molding the Flow of Light. 2nd Edition*, Princeton University Press, Princeton (NJ), 2008.
- [43] Sajeew John, *Strong localization of photons in certain disordered dielectric superlattices*, Phys. Rev. Lett. **58** (1987), no. 23, 2486–2489.
- [44] C. Y. Kao, Stanley Osher, and Eli Yablonovitch, *Maximizing band gaps in two-dimensional photonic crystals by using level set methods*, Appl. Phys. B **81** (2005), 235–244.
- [45] Tosio Kato, *Perturbation Theory for Linear Operators. 2nd Edition*, Grundlehren der mathematischen Wissenschaften, no. 132, Springer, Berlin, 1976.
- [46] Shuichi Kinoshita, *Structural Colors in the Realm of Nature*, World Scientific Publishing, Singapore, 2008.
- [47] Rolf Klein, *Concrete and Abstract Voronoi Diagrams*, Lecture Notes in Computer Science, no. 400, Springer, Berlin, 1989.
- [48] Mark G. Krein, *On certain problems on the maximum and minimum of characteristic values and on the Lyapunov zones of stability*, AMS Transl. Ser. 2 **1** (1955), 163–187.
- [49] Peter Kuchment, *Floquet Theory for Partial Differential Equations*, Birkhäuser, Basel, 1993.
- [50] ———, *Mathematics of photonic crystals*, Mathematical Modeling in Optical Science (Gang Bao, ed.), SIAM, Philadelphia (PA), 2001, pp. 207–272.

- [51] Alexander Kurganov and Eitan Tadmor, *New high-resolution semi-discrete central schemes for Hamilton-Jacobi equations*, J. Comput. Phys. **160** (2000), no. 2, 720–742.
- [52] Gerard Lebourg, *Valeur moyenne pour gradient généralisé*, C. R. Acad. Sci., Paris, Sér. A **281** (1975), 795–797.
- [53] Jerrold E. Marsden and Thomas J.R. Hughes, *Mathematical foundations of elasticity.*, Prentice-Hall, Englewood Cliffs (NJ), 1983.
- [54] Clyde W. Mason, *Structural colors in feathers. II*, J. Phys. Chem. **27** (1923), no. 5, 401–448.
- [55] James C. Maxwell, *A Treatise on Electricity and Magnetism*, Clarendon Press, Oxford, 1873.
- [56] François Murat and Luc Tartar, *H-convergence*, Topics in the Mathematical Modelling of Composite Materials (Andrej Cherkaev and Robert Kohn, eds.), Birkhäuser, Basel, 1997, pp. 21–43.
- [57] Jean-Claude Nédélec, *A new family of mixed finite elements in \mathbb{R}^3* , Numer. Math. **50** (1986), 57–81.
- [58] Stanley Osher and Ronald Fedkiw, *Level Set Methods and Dynamic Implicit Surfaces*, Springer, Berlin, 2003.
- [59] Stanley Osher and Fadil Santosa, *Level set methods for optimization problems involving geometry and constraints. I. Frequency of a two-density inhomogeneous drum*, J. Comput. Phys. **171** (2001), 272–288.
- [60] Dennis W. Prather, Shouyuan Shi, Ahmed Sharkawy, Janusz Murakowski, and Garrett J. Schneider, *Photonic Crystals: Theory, Applications, and Fabrication*, John Wiley & Sons, Hoboken (NJ), 2009.
- [61] Stefan Preble, Michal Lipson, and Hod Lipson, *Two-dimensional photonic crystals designed by evolutionary algorithms*, Appl. Phys. Lett. **86** (2005), 061111.
- [62] Erhard Quaisser, *Geometrie: Einführung, Probleme, Übungen.*, Spektrum Akademischer Verlag, Heidelberg, 1994.
- [63] Lord Rayleigh, *On the dynamical theory of gratings*, Lond. R. S. Proc. (A) **79** (1907), 399–416.

- [64] Michael Reed and Barry Simon, *Methods of Modern Mathematical Physics. IV: Analysis of Operators*, Academic Press, New York, 1978.
- [65] R. Tyrrell Rockafellar, *Convex Analysis*, Princeton University Press, Princeton (NJ), 1970.
- [66] Kazuaki Sakoda, *Optical Properties of Photonic Crystals. 2nd Edition*, Springer, Berlin, 2005.
- [67] Laurent Schwartz, *Radon Measures on Arbitrary Topological Spaces and Cylindrical Measures*, Oxford University Press, Oxford, 1973.
- [68] Edward G. Sewell, *Automatic Generation of Triangulations for Piecewise Polynomial Approximations*, Ph.D. thesis, Purdue University, 1973.
- [69] Ole Sigmund and Jensen Søndergaard, *Systematic design of photonic band-gap materials and structures by topology optimization*, Phil. Trans. R. Soc. Lond. A **361** (2003), 1001–1019.
- [70] Jan Sokolowski and Jean-Paul Zolesio, *Introduction to Shape Optimization: Shape Sensitivity Analysis*, Springer-Verlag, Berlin, 1992.
- [71] Krister Svanberg, *The method of moving asymptotes: A new method for structural optimization*, Int. J. Numer. Methods Eng. **24** (1987), 359–373.
- [72] Luc Tartar, *Compensated compactness and applications to partial differential equations*, Nonlinear Analysis and Mechanics, Heriot-Watt Symposium IV (Robin J. Knops, ed.), Pitman, London, 1979, pp. 136–211.
- [73] Pavel Šolín, Karel Segeth, and Ivo Deležel, *Higher-Order Finite Element Methods*, CRC Press, Boca Raton (FL), 2004.
- [74] Christian Weber, *A local compactness theorem for maxwell's equations*, Math. Methods Appl. Sci. **2** (1980), 12–25.
- [75] Dirk Werner, *Funktionalanalysis. 6. Auflage*, Springer, Berlin, 2007.
- [76] Christian Wieners, *Distributed point objects: A new concept for parallel finite elements*, Domain Decomposition Methods in Science and Engineering (Berlin) (Ralf Kornhuber, Ronald Hoppe, Jacques Périaux, Olivier Pironneau, Olof Widlund, and Jinchao Xu, eds.), Lecture Notes in Computational Science and Engineering, vol. 40, Springer, 2004, pp. 175–183.
- [77] ———, *A geometric data structure for parallel finite elements and the application to multigrid methods with block smoothing*, Comput. Vis. Sci. **13** (2010), 161–175.

- [78] Joseph Wloka, *Partial Differential Equations*, Cambridge University Press, Cambridge, 1987.
- [79] Rob S. Womersley, *Local properties of algorithms for minimizing nonsmooth composite functions*, Math. Program. **32** (1985), 69–89.
- [80] Eli Yablonovitch, *Inhibited spontaneous emission in solid-state physics and electronics*, Phys. Rev. Lett. **58** (1987), no. 20, 2059–2062.