



HAL
open science

Eigenvalue enclosures and applications to the Maxwell operator

Nabile Boussaid, Gabriel Raúl Barrenechea, Lyonell Boulton

► **To cite this version:**

Nabile Boussaid, Gabriel Raúl Barrenechea, Lyonell Boulton. Eigenvalue enclosures and applications to the Maxwell operator. 2013. hal-00837475v1

HAL Id: hal-00837475

<https://hal.science/hal-00837475v1>

Preprint submitted on 22 Jun 2013 (v1), last revised 19 Feb 2014 (v3)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

EIGENVALUE ENCLOSURES AND APPLICATIONS TO THE MAXWELL OPERATOR

GABRIEL R. BARRENECHEA, LYONELL BOULTON, AND NABILE BOUSSAÏD

ABSTRACT. In this work we discuss the numerical estimation of the eigenfrequencies and field phasors of the resonant cavity problem on a bounded region filled with a possibly anisotropic medium. We present a general framework which allows finding lower and upper bounds for the eigenfrequencies, hence providing a computable residual and multiplicity counting. We establish precise rates of convergence of the method in general, and show that this rate is optimal for trial spaces of standard nodal finite elements. In the final part of the paper we include a reproducible computational procedure and then report on various numerical experiments performed on two and three-dimensional benchmark geometries, with and without symmetries.

CONTENTS

1. Introduction	1
2. Approximated local counting functions	4
3. The method of Zimmermann and Mertins	9
4. Convergence and error estimates	12
5. The finite element method for the Maxwell eigenvalue problem	20
6. The numerical strategy in a nutshell	23
7. Computational examples	25
7.1. Convex domains	26
7.2. Non-convex domains	31
Appendix A. Further geometrical properties of $F_{\mathcal{L}}^j(t)$	37
Appendix B. A core Comsol v4.3 LiveLink code	38
Acknowledgements	41
References	41

1. INTRODUCTION

In this work we discuss the numerical computation of eigenvalue enclosures and approximating eigenfunctions for the Maxwell operator on a bounded domain filled with a possibly anisotropic medium. We establish a general framework which allows finding lower and upper bounds for eigenvalues, hence providing a computable residual and multiplicity counting. The origins of this framework can be traced back to the works of Zimmermann and Mertins, [24], and Davies, [16], and it is

Date: June 22, 2013.

Key words and phrases. eigenvalue enclosures, Maxwell equation, spectral pollution, finite element method.

guaranteed to be free from spectral pollution. One of our goals is to determine the precise rate of convergence of the method. This rate turns out to be optimal in the sense of standard interpolation estimates, when a concrete implementation is performed on trial spaces of nodal finite elements. In the final part of the paper we include a reproducible computational procedure and report on the outcomes of several benchmark numerical experiments.

Let $\Omega \subset \mathbb{R}^3$ be a domain which is bounded, open and simply connected. Everywhere below we assume that $\partial\Omega$, the boundary of Ω , is sufficiently regular (in a sense made precise in Section 5) and denote by \mathbf{n} its outer normal vector. Consider the Maxwell eigenvalue problem

$$(1) \quad \begin{cases} \operatorname{curl} \mathbf{E} = i\omega\mu\mathbf{H} & \text{in } \Omega \\ \operatorname{curl} \mathbf{H} = -i\omega\epsilon\mathbf{E} & \text{in } \Omega \\ \mathbf{E} \times \mathbf{n} = 0 & \text{on } \partial\Omega \end{cases}$$

where the angular frequencies $\omega \in \mathbb{R}$ are to be determined alongside with the electric and magnetic field phasors $(\mathbf{E}, \mathbf{H}) \neq 0$. Here and elsewhere the electric permittivity and the magnetic permeability, ϵ and μ respectively, are positive and such that

$$(2) \quad \epsilon, \frac{1}{\epsilon}, \mu, \frac{1}{\mu} \in L^\infty(\Omega).$$

The electromagnetic oscillations in a resonator are described by equation (1) restricted to the solenoidal subspace characterized by Gauss's law,

$$(3) \quad \operatorname{div}(\epsilon\mathbf{E}) = 0 = \operatorname{div}(\mu\mathbf{H}).$$

The orthogonal complement of the latter in a suitably weighted L^2 space (see [7] and references therein) is the gradient space, which has infinite dimension and it lies in the kernel of the self-adjoint operator \mathcal{M} associated to (1). In turns, this means that the non-zero spectrum and non-zero eigenspace of (1), with or without the ansatz (3), are identical. Below we propose computing the non-zero angular frequencies and field phasors of the resonator by means of the unrestricted problem (1).

The operator \mathcal{M} does not have a compact resolvent and it is strongly indefinite. The self-adjoint operator associated to (1)-(3) has a compact resolvent but it is still strongly indefinite. By considering the square of \mathcal{M} restricted to the solenoidal space, one obtains a positive definite eigenvalue problem (involving the bi-curl, for example, if the medium is isotropic) which can in principle be discretized via the Galerkin method. One serious drawback of this idea for practical computations is the fact that the standard finite element spaces are not solenoidal. Usually, spurious modes associated to the infinite-dimensional kernel appear and give rise to spectral pollution. This has been well documented to be a manifested problem whenever the underlying mesh is unstructured, [2] and references therein.

In order to avoid spectral pollution, various ingenious methods have been considered in the past. One possible approach [10], is to enhance the divergence of the electric field in a fractional order negative Sobolev norm. In [12] nodal elements are employed in conjunction with a least squares formulation of a weak form re-writting of (1)-(3). The condition (3) can also be incorporated into (1) by means of a Lagrange multiplier and then employ continuous finite element spaces of a Taylor-Hood type, [13]. A further possible approach, [9, 8], is to re-write the

spectral problem associated to \mathcal{M}^2 in a mixed form and then consider edge finite elements. The latter is widely regarded as the most effective way to avoid spurious modes for the resonant cavity problem. Moreover, it turns out to be linked to deep results on the rigorous treatment of finite elements and it is at the core of elegant geometrical ideas, [2].

Unfortunately, as far as we are aware, all these computational techniques available for the Maxwell problem currently exhibit two main limitations.

- a) They are not certified. To be precise, computed eigenvalues are not necessarily guaranteed one-sided bounds of the exact eigenvalues in general, despite of the possible convergence of the method.
- b) Detecting the multiplicity of an eigenvalue or the presence of a cluster of eigenvalues is extremely difficult.

Our goal below is to propose an alternative approach for computing the eigenvalues of (1) which addresses these limitations. The general strategy is based on the extensions of the Temple-Lehmann-Goerisch method [21] developed by Zimmermann and Mertins, [24]. We show that the procedure is robust in the sense that any standard class of finite elements, including the ones based on nodal degrees of freedom, can be employed to perform computations which are certified up to machine precision. In recent years, this method has been successfully implemented in the context of the radially reduced magnetohydrodynamics operator [24, 11], the Helmholtz equation [5] and the calculation of sloshing frequencies in the left definite case [4].

The method of Zimmermann and Mertins is closely linked to another pollution-free technique for eigenvalue computation which is based on a notion of approximated spectral distance. This other method has been examined by Davies in [16, 17] and later by Davies and Plum in [18], but it is yet to be tested properly on models of dimension other than one. Below we develop further the arguments presented in [18, Section 6], in order to demonstrate that these two approaches are equivalent.

In Section 2 we extend various results of [18]. Notably, we include in our current analysis the determination of multiplicities for eigenvalues in Proposition 1 and a description of how eigenfunctions are approximated in Proposition 3. We introduce the method of Zimmermann and Mertins in Section 3. We derive the latter in a self-contained manner independently from the work [24]. See Theorem 6 and Corollary 7.

In Section 4 we examine convergence and residuals. The main statement of this section is Theorem 12, where we determine a general convergence estimate with explicit residuals for a finite group of contiguous eigenvalues. This theorem is employed in Section 5, where we establish concrete approximation rates for the pollution-free numerical solution of (1) by means of nodal finite elements. Theorem 14 collects the main contribution in this respect. We show that the rate of convergence achieved by the current method is optimal for these trial spaces. Here we have chosen the most widely available class of finite elements, in order to illustrate our findings in a concrete accessible manner. However we should remark that the technique described in sections 2 and 3 is formulated in general, and it can also be implemented on other classes of basis functions.

The final part of the paper has a more computational character. A concrete numerical strategy is specified in the Procedure 1 of Section 6. According to

Lemma 16, this strategy is convergent in a suitable regime for finite elements. Section 7, on the other hand, is devoted to various benchmark experiments on specific models which demonstrate the applicability of the proposed technique.

2. APPROXIMATED LOCAL COUNTING FUNCTIONS

This section is devoted to notions of approximated spectral distance and approximated local counting function for self-adjoint operators. We follow closely the framework established in [16, 17, 18]. These notions and their properties will lead in the next section to the formulation of a method for eigenvalue computation which has been examined in [24] and subsequent works [5, 4]. Various results in all these references can be recovered from the unified approach presented below.

Let $A : D(A) \rightarrow \mathcal{H}$ be a self-adjoint operator acting on a Hilbert space \mathcal{H} . Below we decompose the spectrum of A in the usual fashion, as the union of discrete and essential spectrum, $\sigma(A) = \sigma_{\text{disc}}(A) \cup \sigma_{\text{ess}}(A)$. Let J be any Borel subset of \mathbb{R} . The spectral projector associated to A is denoted by $\mathbb{1}_J(A) = \int_J dE_\lambda$. Hence $\text{Tr } \mathbb{1}_J(A) = \dim \mathbb{1}_J(A)\mathcal{H}$. We write $\mathcal{E}_J(A) = \bigoplus_{\lambda \in J} \ker(A - \lambda)$ with the convention $\mathcal{E}_\lambda(A) = \mathcal{E}_{\{\lambda\}}(A)$. Generally $\mathcal{E}_J(A) \subseteq \mathbb{1}_J(A)\mathcal{H}$, however there is no reason for these two subspaces to be equal.

Let $t \in \mathbb{R}$. Let $q_t : D(A) \times D(A) \rightarrow \mathbb{C}$ be the closed bilinear form

$$(4) \quad q_t(u, w) = \langle (A - t)u, (A - t)w \rangle \quad \forall u, w \in D(A).$$

For any $u \in D(A)$ we will constantly make use of the following t -dependant seminorm, which is a norm if t is not an eigenvalue,

$$(5) \quad |u|_t = q_t(u, u)^{1/2} = \|(A - t)u\|.$$

By virtue of the min-max principle, q_t characterizes the spectrum which lies near the origin of the positive operator $(A - t)^2$. In turn, this gives rise to a notion of local counting function at t for the spectrum of A .

Let

$$\mathfrak{d}_j(t) = \inf_{\substack{\dim V=j \\ V \subset D(A)}} \sup_{u \in V} \frac{|u|_t}{\|u\|}$$

so that $0 \leq \mathfrak{d}_j(t) \leq \mathfrak{d}_k(t)$ for $j < k$. Then $\mathfrak{d}_1(t)$ is the Hausdorff distance from t to $\sigma(A)$,

$$(6) \quad \mathfrak{d}_1(t) = \min\{\lambda \in \sigma(A) : |\lambda - t|\} = \inf_{u \in D(A)} \frac{|u|_t}{\|u\|}.$$

Similarly $\mathfrak{d}_j(t)$ are the distances from t to the j th nearest point in $\sigma(A)$ counting multiplicity in a generalized sense. That is, stopping when the essential spectrum is reached. Moreover

$$\mathfrak{d}_j(t) = \mathfrak{d}_{j-1}(t) \iff \begin{cases} \text{either} & \dim \mathcal{E}_{[t - \mathfrak{d}_{j-1}(t), t + \mathfrak{d}_{j-1}(t)]}(A) > j - 1 \\ \text{or} & t + \mathfrak{d}_{j-1}(t) \in \sigma_{\text{ess}}(A) \\ \text{or} & t - \mathfrak{d}_{j-1}(t) \in \sigma_{\text{ess}}(A). \end{cases}$$

Without further mention, below we will always count spectral points of A relative to t , regarding multiplicities in this generalized sense.

We now show how to extract certified information about $\sigma(A)$ in the vicinity of t from the action of A onto finite-dimensional trial subspaces $\mathcal{L} \subset \mathcal{D}(A)$, see [16, Section 3]. For $j \leq n = \dim \mathcal{L}$, let

$$(7) \quad F_{\mathcal{L}}^j(t) = \min_{\substack{\dim V=j \\ V \subset \mathcal{L}}} \max_{u \in V} \frac{|u|_t}{\|u\|}.$$

Then $0 \leq F_{\mathcal{L}}^1(t) \leq \dots \leq F_{\mathcal{L}}^n(t)$ and $F_{\mathcal{L}}^j(t) \geq \mathfrak{d}_j(t)$ for all $j = 1, 2, \dots, n$. Since $[t - \mathfrak{d}_j(t), t + \mathfrak{d}_j(t)] \subseteq [t - F_{\mathcal{L}}^j(t), t + F_{\mathcal{L}}^j(t)]$, there are at least j spectral points of A in the segment $[t - F_{\mathcal{L}}^j(t), t + F_{\mathcal{L}}^j(t)]$ including, possibly, the essential spectrum. That is

$$(8) \quad \text{Tr } \mathbb{1}_{[t - F_{\mathcal{L}}^j(t), t + F_{\mathcal{L}}^j(t)]}(A) \geq j \quad \forall j = 1, \dots, n.$$

Hence $F_{\mathcal{L}}^j(t)$ is an approximated local counting function for $\sigma(A)$.

As a consequence of the triangle inequality, $F_{\mathcal{L}}^j$ is a Lipschitz continuous function such that

$$(9) \quad |F_{\mathcal{L}}^j(t) - F_{\mathcal{L}}^j(s)| \leq |t - s| \quad \forall s, t \in \mathbb{R} \quad \text{and} \quad j = 1, \dots, n.$$

Moreover, $F_{\mathcal{L}}^j(t)$ is the j th smallest eigenvalue μ of the non-negative weak problem:

$$(10) \quad \text{find } (\mu, u) \in [0, \infty) \times \mathcal{L} \setminus \{0\} \quad \text{such that} \quad q_t(u, v) = \mu^2 \langle u, v \rangle \quad \forall v \in \mathcal{L}.$$

Hence

$$(11) \quad F_{\mathcal{L}}^j(t) = \max_{\substack{\dim V=j-1 \\ V \subset \mathcal{L}}} \min_{u \in \mathcal{L} \ominus V} \frac{|u|_t}{\|u\|} = \max_{\substack{\dim V=j-1 \\ V \subset \mathcal{H}}} \min_{u \in \mathcal{L} \ominus V} \frac{|u|_t}{\|u\|}.$$

We now show how to detect the spectrum of A to the left/right of t by means of $F_{\mathcal{L}}^j$ in an optimal setting. This turns out to be a crucial ingredient in the formulation of the strategy proposed in [16, 17, 18]. The following notation simplifies various statements below. Let

$$\begin{aligned} \mathfrak{n}_j^-(t) &= \sup\{s < t : \text{Tr } \mathbb{1}_{(s,t]}(A) \geq j\} \quad \text{and} \\ \mathfrak{n}_j^+(t) &= \inf\{s > t : \text{Tr } \mathbb{1}_{[t,s)}(A) \geq j\}. \end{aligned}$$

Then $\mathfrak{n}_j^{\mp}(t)$ is the j th point in $\sigma(A)$ to the left(-)/right(+) of t counting multiplicities. Here $t \in \sigma(A)$ is allowed and neither t nor $\mathfrak{n}_1^{\mp}(t)$ have to be isolated from the rest of $\sigma(A)$. Note that $\mathfrak{n}_j^-(t) = -\infty$ for $\text{Tr } \mathbb{1}_{(-\infty,t]}(A) < j$ and $\mathfrak{n}_j^+(t) = +\infty$ for $\text{Tr } \mathbb{1}_{[t,+\infty)}(A) < j$. Without further mention, all statements below regarding bounds on $\mathfrak{n}_j^{\mp}(t)$ will be void (hence redundant) in either of these two cases.

Proposition 1. *Let $t^- < t < t^+$. Then*

$$(12) \quad \begin{aligned} F_{\mathcal{L}}^j(t^-) \leq t - t^- &\quad \Rightarrow \quad t^- - F_{\mathcal{L}}^j(t^-) \leq \mathfrak{n}_j^-(t) \\ F_{\mathcal{L}}^j(t^+) \leq t^+ - t &\quad \Rightarrow \quad t^+ + F_{\mathcal{L}}^j(t^+) \geq \mathfrak{n}_j^+(t). \end{aligned}$$

Moreover, let $t_1^- < t_2^- < t < t_2^+ < t_1^+$. Then

$$(13) \quad \begin{aligned} F_{\mathcal{L}}^j(t_i^-) \leq t - t_i^- \text{ for } i = 1, 2 &\quad \Rightarrow \quad t_1^- - F_{\mathcal{L}}^j(t_1^-) \leq t_2^- - F_{\mathcal{L}}^j(t_2^-) \leq \mathfrak{n}_j^-(t) \\ F_{\mathcal{L}}^j(t_i^+) \leq t_i^+ - t \text{ for } i = 1, 2 &\quad \Rightarrow \quad t_1^+ + F_{\mathcal{L}}^j(t_1^+) \geq t_2^+ + F_{\mathcal{L}}^j(t_2^+) \geq \mathfrak{n}_j^+(t). \end{aligned}$$

Proof. We firstly show (12). Suppose that $t \geq F_{\mathcal{L}}^j(t^-) + t^-$. Then

$$\mathrm{Tr} \mathbb{1}_{[t^- - F_{\mathcal{L}}^j(t^-), t]}(A) \geq j.$$

Since $\mathbf{n}_j^-(t) \leq \dots \leq \mathbf{n}_1^-(t)$ are the only spectral points in the segment $[\mathbf{n}_j^-(t), t]$, then necessarily

$$\mathbf{n}_j^-(t) \in [t^- - F_{\mathcal{L}}^j(t^-), t].$$

The bottom of (12) is shown in a similar fashion.

The second statement follows by observing that the maps $t \mapsto t \pm F_{\mathcal{L}}^j(t)$ are monotonically increasing as a consequence of (9). \square

The structure of the trial subspace \mathcal{L} determines the existence of t^\pm satisfying the hypothesis in (12). If we expect to detect $\sigma(A)$ at both sides of t , a necessary requirement on \mathcal{L} should certainly be the condition

$$(14) \quad \min_{u \in \mathcal{L}} \frac{\langle Au, u \rangle}{\langle u, u \rangle} < t < \max_{u \in \mathcal{L}} \frac{\langle Au, u \rangle}{\langle u, u \rangle}.$$

By virtue of lemmas 4 and 5 below, for $j = 1$, the left hand side inequality of (14) implies the existence of t^- and the right hand side inequality implies the existence of t^+ , respectively.

Remark 1. From Proposition 1 it follows that optimal lower bounds for $\mathbf{n}_j^-(t)$ are achieved by finding $\hat{t}_j^- \leq t$, the closer point to t , such that $F_{\mathcal{L}}^j(\hat{t}_j^-) = t - \hat{t}_j^-$. Indeed, by virtue of (13), $t^- - F_{\mathcal{L}}^j(t^-) \leq \hat{t}_j^- - F_{\mathcal{L}}^j(\hat{t}_j^-) \leq \mathbf{n}_j^-(t)$ for any other t^- as in (12). Similarly, optimal upper bounds for $\mathbf{n}_j^+(t)$ are found by analogous means. This observation will play a crucial role in Section 3.

We now determine further geometrical properties of $F_{\mathcal{L}}^1$ and its connection to the spectral distance. Let the Hausdorff distances from $t \in \mathbb{R}$ to $\sigma(A) \setminus (-\infty, t]$ and $\sigma(A) \setminus [t, \infty)$, respectively, be given by

$$(15) \quad \begin{aligned} \delta^+(t) &= \inf\{\mu - t : \mu \in \sigma(A), \mu > t\} \quad \text{and} \\ \delta^-(t) &= \inf\{t - \mu : \mu \in \sigma(A), \mu < t\}. \end{aligned}$$

In general, $t - \mathbf{n}_1^-(t) \leq \delta^-(t)$ and $\mathbf{n}_1^+(t) - t \leq \delta^+(t)$. In fact, $|\mathbf{n}_1^\pm(t) - t| = \delta^\pm(t)$ for $t \notin \sigma(A)$. However, these relations can be strict whenever $t \in \sigma(A)$. Indeed, $\mathbf{n}_1^+(t) - t = \delta^+(t)$ iff there exists a decreasing sequence $t_n^+ \in \sigma(A)$ such that $t_n^+ \downarrow t$, whereas $t - \mathbf{n}_1^-(t) = \delta^-(t)$ iff there exists an increasing sequence $t_n^- \in \sigma(A)$ such that $t_n^- \uparrow t$.

An emphasis in distinguishing $|\mathbf{n}_1^\pm(t) - t|$ from $\delta^\pm(t)$ seems unnecessary at this stage. However, this distinction in the notation will be justified later on. Without further mention below we write $\delta^\pm(t) = \pm\infty$ to indicate that either of the sets on the right side of (15) is empty.

Let $\lambda \in \sigma(A)$ be an isolated point. If there exists a non-vanishing $u \in \mathcal{L} \cap \mathcal{E}_\lambda(A)$, then

$$\frac{|u|_s}{\|u\|} = |\lambda - s| = \mathfrak{d}_1(s) \quad \forall s \in \left[\lambda - \frac{\delta^-(\lambda)}{2}, \lambda + \frac{\delta^+(\lambda)}{2} \right].$$

According to the convergence analysis carried out in Section 4, the smaller the angle between \mathcal{L} and the spectral subspace $\mathcal{E}_\lambda(A)$, the closer the $F_{\mathcal{L}}^1(t)$ is to $\mathfrak{d}_1(t)$ for $t \in \left(\lambda - \frac{\delta^-(\lambda)}{2}, \lambda + \frac{\delta^+(\lambda)}{2} \right)$. The special case of this angle being zero is described by the following lemma.

Lemma 2. *For $\lambda \in \sigma(A)$ isolated from the rest of the spectrum, the following statements are equivalent.*

- a) *There exists a minimizer $u \in \mathcal{L}$ of the right side of (7) for $j = 1$, such that $|u|_t = \mathfrak{d}_1(t)$ for a single $t \in (\lambda - \frac{\delta^-(\lambda)}{2}, \lambda + \frac{\delta^+(\lambda)}{2})$,*
- b) *$F_{\mathcal{L}}^1(t) = \mathfrak{d}_1(t)$ for a single $t \in (\lambda - \frac{\delta^-(\lambda)}{2}, \lambda + \frac{\delta^+(\lambda)}{2})$,*
- c) *$F_{\mathcal{L}}^1(s) = \mathfrak{d}_1(s)$ for all $s \in [\lambda - \frac{\delta^-(\lambda)}{2}, \lambda + \frac{\delta^+(\lambda)}{2}]$,*
- d) *$\mathcal{L} \cap \mathcal{E}_\lambda(A) \neq \{0\}$.*

Proof. Since \mathcal{L} is finite-dimensional, a) and b) are equivalent by the definitions of $\mathfrak{d}_1(t)$, $F_{\mathcal{L}}^1(t)$ and q_t . From the paragraph above the statement of the lemma it is clear that d) \Rightarrow c) \Rightarrow b). Since $|u|_t/\|u\|$ is the square root of the Rayleigh quotient associated to the operator $(A - t)^2$, the fact that λ is isolated combined with the Rayleigh-Ritz principle, gives the implication a) \Rightarrow d). \square

As there can be a mixing of eigenspaces, it is not possible to replace b) in this lemma by an analogous statement including $t = \lambda \pm \frac{\delta^\pm(\lambda)}{2}$. If $\lambda' = \lambda + \delta^+(\lambda)$ is an eigenvalue, for example, then $F_{\mathcal{L}}^1\left(\frac{\lambda + \lambda'}{2}\right) = \mathfrak{d}_1\left(\frac{\lambda + \lambda'}{2}\right)$ ensures that \mathcal{L} contains elements of $\mathcal{E}_\lambda(A) \oplus \mathcal{E}_{\lambda'}(A)$. However it is not guaranteed to be orthogonal to either of these two subspaces. See the Appendix A for similar results in the case $j > 1$.

We conclude this section by examining extensions of the implications b) \Rightarrow d) of Lemma 2 into a more general context. In combination with the results of Section 3, the next proposition shows how to obtain certified information about spectral subspaces. Some of its practical implications will be discussed later on in Section 7.

Here and below $\{u_j^t\}_{j=1}^n \subset \mathcal{L}$ will denote an orthonormal family of eigenfunctions associated to the eigenvalues $\mu = F_{\mathcal{L}}^j(t)$ of the weak problem (10). In a suitable asymptotic regime for \mathcal{L} , the angle between these eigenfunctions and the spectral subspaces of $|A - t|$ in the vicinity of the origin is controlled by a residual which is as small as $\mathcal{O}\left(\sqrt{F_{\mathcal{L}}^j(t) - \mathfrak{d}_j(t)}\right)$ for $F_{\mathcal{L}}^j(t) - \mathfrak{d}_j(t) \rightarrow 0$.

Assumption 1. *Unless otherwise specified, from now on we will always fix the parameter $m \leq n = \dim \mathcal{L}$ and suppose that*

$$(16) \quad [t - \mathfrak{d}_m(t), t + \mathfrak{d}_m(t)] \cap \sigma(A) \subseteq \sigma_{\text{disc}}(A).$$

Set

$$\delta_j(t) = \text{dist} \left[t, \sigma(A) \setminus \{t \pm \mathfrak{d}_k(t)\}_{k=1}^j \right].$$

By virtue of (16), $\delta_j(t) > \mathfrak{d}_j(t)$ for all $j \leq m$.

Remark 2. If $t = \frac{n_i^-(t) + n_j^+(t)}{2}$ for a given j , the vectors ϕ_j^t introduced in Proposition 3 and invoked subsequently, might not be eigenvectors of A despite of the fact that $|A - t|\phi_j^t = \mathfrak{d}_j(t)\phi_j^t$. However, in any other circumstance ϕ_j^t are eigenvectors of A .

Proposition 3. *Let $t \in \mathbb{R}$ and $j \in \{1, \dots, m\}$. Assume that $F_{\mathcal{L}}^j(t) - \mathfrak{d}_j(t)$ is small enough so that $0 < \varepsilon_j < 1$ holds true for the residuals constructed inductively as*

follows,

$$\varepsilon_1 = \sqrt{\frac{F_{\mathcal{L}}^1(t)^2 - \mathfrak{d}_1(t)^2}{\delta_1(t)^2 - \mathfrak{d}_1(t)^2}}$$

$$\varepsilon_j = \sqrt{\frac{F_{\mathcal{L}}^j(t)^2 - \mathfrak{d}_j(t)^2}{\delta_j(t)^2 - \mathfrak{d}_j(t)^2} + \sum_{k=1}^{j-1} \frac{\varepsilon_k^2}{1 - \varepsilon_k^2} \left(1 + \frac{\mathfrak{d}_j(t)^2 - \mathfrak{d}_k(t)^2}{\delta_j(t)^2 - \mathfrak{d}_j(t)^2}\right)}.$$

Then, there exists an orthonormal basis $\{\phi_j^t\}_{j=1}^m$ of $\mathcal{E}_{[t-\mathfrak{d}_m(t), t+\mathfrak{d}_m(t)]}(A)$ such that $\phi_j^t \in \mathcal{E}_{\{t-\mathfrak{d}_j(t), t+\mathfrak{d}_j(t)\}}(A)$,

$$(17) \quad \|u_j^t - \langle u_j^t, \phi_j^t \rangle \phi_j^t\| \leq \varepsilon_j \quad \text{and}$$

$$(18) \quad |u_j^t - \langle u_j^t, \phi_j^t \rangle \phi_j^t|_t \leq \sqrt{F_{\mathcal{L}}^j(t)^2 - \mathfrak{d}_j(t)^2 + \mathfrak{d}_j(t)^2 \varepsilon_j^2}.$$

Proof. As it is clear from the context, in this proof we suppress the index t on top of any vector. We write $\Pi_{\mathcal{S}}$ to denote the orthogonal projection onto the subspace \mathcal{S} with respect to the inner product $\langle \cdot, \cdot \rangle$.

Let us first consider the case $j = 1$. Let $\mathcal{S}_1 = \mathcal{E}_{\{t-\mathfrak{d}_1(t), t+\mathfrak{d}_1(t)\}}(A)$, and decompose $u_1 = \Pi_{\mathcal{S}_1} u_1 + u_1^\perp$ where $u_1^\perp \perp \mathcal{S}_1$. Since A is self-adjoint,

$$(19) \quad F_{\mathcal{L}}^1(t)^2 = \|(A-t)u_1\|^2 = \mathfrak{d}_1(t)^2 \|\Pi_{\mathcal{S}_1} u_1\|^2 + \|(A-t)u_1^\perp\|^2.$$

Hence

$$F_{\mathcal{L}}^1(t)^2 \geq \mathfrak{d}_1(t)^2 (1 - \|u_1^\perp\|^2) + \delta_1(t)^2 \|u_1^\perp\|^2.$$

Since $\delta_1(t) > \mathfrak{d}_1(t)$, clearing from this identity $\|u_1^\perp\|^2$ yields $\|u_1^\perp\| \leq \varepsilon_1$. Hence $\|\Pi_{\mathcal{S}_1} u_1\|^2 \geq 1 - \varepsilon_1^2 > 0$. Let

$$\phi_1 = \frac{1}{\|\Pi_{\mathcal{S}_1} u_1\|} \Pi_{\mathcal{S}_1} u_1$$

so that $\|\Pi_{\mathcal{S}_1} u_1\| = |\langle u_1, \phi_1 \rangle|$. Then (17) holds immediately and (18) is achieved by clearing $\|(A-t)u_1^\perp\|^2$ from (19).

We define the needed basis, and show (17) and (18), for j up to m inductively as follows. Set

$$\phi_j = \frac{1}{\|\Pi_{\mathcal{S}_j} u_j\|} \Pi_{\mathcal{S}_j} u_j$$

where $\mathcal{S}_j = \mathcal{E}_{\{t-\mathfrak{d}_j(t), t+\mathfrak{d}_j(t)\}}(A) \ominus \text{Span}\{\phi_l\}_{l=1}^{j-1}$ and $\Pi_{\mathcal{S}_j} u_j \neq 0$, all this for $1 \leq j \leq k-1$. Assume that (17) and (18) hold true for j up to $k-1$. Define $\mathcal{S}_k = \mathcal{E}_{\{t-\mathfrak{d}_k(t), t+\mathfrak{d}_k(t)\}}(A) \ominus \text{Span}\{\phi_l\}_{l=1}^{k-1}$. We first show that $\Pi_{\mathcal{S}_k} u_k \neq 0$, and so we can define

$$(20) \quad \phi_k = \frac{1}{\|\Pi_{\mathcal{S}_k} u_k\|} \Pi_{\mathcal{S}_k} u_k$$

ensuring $\phi_k \perp \text{Span}\{\phi_l\}_{l=1}^{k-1}$. After that we verify the validity of (17) and (18) for $j = k$.

Decompose

$$u_k = \Pi_{\mathcal{S}_k} u_k + \sum_{l=k-1}^1 \langle u_k, \phi_l \rangle \phi_l + u_k^\perp$$

where $u_k^\perp \perp \text{Span}\{\phi_l\}_{l=1}^{k-1} \oplus \mathcal{S}_k$. Then

$$\begin{aligned} F_{\mathcal{L}}^k(t)^2 &= \mathfrak{d}_k(t)^2 \|\Pi_{\mathcal{S}_k} u_k\|^2 + \sum_{l=k-1}^1 \mathfrak{d}_l(t)^2 |\langle u_k, \phi_l \rangle|^2 + \|(A-t)u_k^\perp\|^2 \\ &\geq \mathfrak{d}_k(t)^2 \|\Pi_{\mathcal{S}_k} u_k\|^2 + \sum_{l=k-1}^1 \mathfrak{d}_l(t)^2 |\langle u_k, \phi_l \rangle|^2 + \delta_k(t)^2 \|u_k^\perp\|^2 \\ &= \mathfrak{d}_k(t)^2 (1 - \|u_k^\perp\|^2) + \sum_{l=k-1}^1 (\mathfrak{d}_l(t)^2 - \mathfrak{d}_k(t)^2) |\langle u_k, \phi_l \rangle|^2 + \delta_k(t)^2 \|u_k^\perp\|^2. \end{aligned}$$

The conclusion (17) up to $k-1$, implies $|\langle u_l, \phi_l \rangle|^2 \geq 1 - \varepsilon_l^2$ for $l = 1, \dots, k-1$. Since $\langle u_k, u_l \rangle = 0$ for $l \neq k$,

$$|\langle u_l, \phi_l \rangle| |\langle u_k, \phi_l \rangle| = |\langle u_k, u_l - \langle u_l, \phi_l \rangle \phi_l \rangle|.$$

Then, the Cauchy-Schwarz inequality alongside with (17) yield

$$(21) \quad |\langle u_k, \phi_l \rangle|^2 \leq \frac{\varepsilon_l^2}{1 - \varepsilon_l^2}.$$

Hence, since $\mathfrak{d}_l(t) \leq \mathfrak{d}_k(t)$,

$$F_{\mathcal{L}}^k(t)^2 \geq \mathfrak{d}_k(t)^2 + \sum_{l=k-1}^1 (\mathfrak{d}_l(t)^2 - \mathfrak{d}_k(t)^2) \frac{\varepsilon_l^2}{1 - \varepsilon_l^2} + (\delta_k(t)^2 - \mathfrak{d}_k(t)^2) \|u_k^\perp\|^2.$$

Clearing $\|u_k^\perp\|^2$ from this inequality and combining with the validity of (21) and (17) up to $k-1$, yields $\Pi_{\mathcal{S}_k} u_k \neq 0$.

Let ϕ_k be as in (20). Then (17) is guaranteed for $j = k$. On the other hand, (17) up to $j = k$, (21) and the identity

$$F_{\mathcal{L}}^k(t)^2 = \mathfrak{d}_k(t)^2 |\langle u_k, \phi_k \rangle|^2 + \|(A-t)(u_k - \langle u_k, \phi_k \rangle \phi_k)\|^2,$$

yield (18) up to $j = k$. \square

The main result of this section is Proposition 1, which is central to the hierarchical method for finding eigenvalue inclusions examined a few years ago in [16, 17]. For fixed \mathcal{L} this method leads to bounds for eigenvalues which are far sharper than those obtained from the obvious idea of estimating local minima of $F_{\mathcal{L}}^1(t)$. It was later shown [18] that this hierarchical method is equivalent to another method established in [24], which extends to the indefinite case the classical Temple-Lehmann-Goerisch inequality. From an abstract perspective, Proposition 1 provides an intuitive insight on the mechanism for determining complementary bounds for eigenvalues (in the left definite case, for example). Even though the method proposed in [16, 17, 18] is yet to be explored more systematically in the practical setting, in most circumstances the following technique appears to be easier to implement.

3. THE METHOD OF ZIMMERMANN AND MERTINS

Let $t \in \mathbb{R}$ and $\mathcal{L} \subset D(A)$ be a specified trial subspace as above. Recall that q_t is given by (4). Let $l_t : D(A) \times D(A) \rightarrow \mathbb{C}$ be the (generally not closed) bilinear form associated to $(A-t)$,

$$l_t(u, w) = \langle (A-t)u, w \rangle \quad \forall u, w \in D(A).$$

Our next purpose is to characterize the optimal parameters t^\pm in Proposition 1 as described in Remark 1 by means of the following weak eigenvalue problem,

$$(Z_t^{\mathcal{L}}) \quad \begin{aligned} &\text{find } u \in \mathcal{L} \setminus \{0\} \text{ and } \tau \in \mathbb{R} \text{ such that} \\ &\tau q_t(u, v) = l_t(u, v) \quad \forall v \in \mathcal{L}. \end{aligned}$$

This problem is central to the method of eigenvalue bounds calculation examined in [24] and it will be at the core of the numerical strategy presented in Section 6.

Let

$$\tau_1^-(t) \leq \dots \leq \tau_{n^-}^-(t) < 0 \quad \text{and} \quad 0 < \tau_{n^+}^+(t) \leq \dots \leq \tau_1^+(t),$$

be the negative and positive eigenvalues of $(Z_t^{\mathcal{L}})$ respectively. Here and below $n^\mp(t)$ is the number of these negative and positive eigenvalues, which are both locally constant in t . Below we will denote eigenfunctions associated with $\tau_j^\mp(t)$ by $u_j^\mp(t)$.

The hypotheses (14) ensure the existence of $\tau_1^\mp(t)$. A more concrete connection with the framework of Section 2 is made precise in the following lemma. Its proof is straightforward, hence omitted.

Lemma 4. *In the following lists, the conditions stated are equivalent.*

$$\begin{array}{ll} a^-) & F_{\mathcal{L}}^1(s) > t - s \text{ for all } s < t \\ b^-) & \frac{\langle Au, u \rangle}{\langle u, u \rangle} > t \text{ for all } u \in \mathcal{L} \\ c^-) & \text{all the eigenvalues of } (Z_t^{\mathcal{L}}) \\ & \text{are positive,} \end{array} \quad \begin{array}{ll} a^+) & F_{\mathcal{L}}^1(s) < s - t \text{ for all } s > t \\ b^+) & \frac{\langle Au, u \rangle}{\langle u, u \rangle} < t \text{ for all } u \in \mathcal{L} \\ c^+) & \text{all the eigenvalues of } (Z_t^{\mathcal{L}}) \\ & \text{are negative.} \end{array}$$

Remark 3. Let $\mathcal{L} = \text{Span}\{b_j\}_{j=1}^n$. The matrix $[q_t(b_j, b_k)]_{j,k=1}^n$ is singular if and only if $\mathcal{E}_t(A) \cap \mathcal{L} \neq \{0\}$. On the other hand, the kernel of $(Z_t^{\mathcal{L}})$ might be non-empty. If $n_0(t)$ is the dimension of this kernel and $n_\infty(t) = \dim(\mathcal{E}_t(A) \cap \mathcal{L})$, then $n = n_\infty(t) + n_0(t) + n^-(t) + n^+(t)$.

Assumption 2. *Note that $n_\infty(t) \geq 1$ if and only if $F_{\mathcal{L}}^j(t) = 0$ for $j = 1, \dots, n_\infty(t)$. In this case the conclusions of Lemma 5 and Theorem 6 below become void. In order to write our statements in a more transparent fashion, without further mention from now on we will suppose that*

$$(22) \quad \mathcal{L} \cap \mathcal{E}_t(A) = \{0\}.$$

By virtue of the next three results, finding the eigenvalues of $(Z_t^{\mathcal{L}})$ is equivalent to finding $s = \hat{t}_j^\pm \in \mathbb{R}$ such that

$$(23) \quad t - s = \mp F_{\mathcal{L}}^j(s),$$

and in this case $\hat{t}_j^\pm = t + \frac{1}{2\tau_j^\pm(t)}$. It then follows from Remark 1 that $(Z_t^{\mathcal{L}})$ encodes information about the optimal bounds for the spectrum around t , achievable by (13) in Proposition 1.

We begin with the case $j = 1$, see [18, Theorem 11].

Lemma 5. *Let $t \in \mathbb{R}$.*

- (−) *The smallest eigenvalue $\tau = \tau_1^-(t)$ of $(Z_t^{\mathcal{L}})$ is negative if and only if there exists $s < t$ such that $F_{\mathcal{L}}^1(s) = t - s$. In this case $s = t + \frac{1}{2\tau_1^-(t)}$ and*

$$F_{\mathcal{L}}^1(s) = -\frac{1}{2\tau_1^-(t)} = \frac{|u_1^-(t)|_s}{\|u_1^-(t)\|}$$

for $u = u_1^-(t) \in \mathcal{L}$ the corresponding eigenvector.

(+) The largest eigenvalue $\tau = \tau_1^+(t)$ of $(Z_t^{\mathcal{L}})$ is positive if and only if there exists $s > t$ such that $F_{\mathcal{L}}^1(s) = s - t$. In this case $s = t + \frac{1}{2\tau_1^+(t)}$ and

$$F_{\mathcal{L}}^1(s) = \frac{1}{2\tau_1^+(t)} = \frac{|u_1^+(t)|_s}{\|u_1^+(t)\|}$$

for $u = u_1^+(t) \in \mathcal{L}$ the corresponding eigenvector.

Proof. We only show (–), as the proof of (+) is similar. For all $u \in \mathcal{L}$ and $s \in \mathbb{R}$,

$$q_s(u, u) - F_{\mathcal{L}}^1(s)^2 \langle u, u \rangle = q_t(u, u) + 2(t - s)l_t(u, u) + ((t - s)^2 - F_{\mathcal{L}}^1(s)^2) \langle u, u \rangle.$$

Suppose that $F_{\mathcal{L}}^1(s) = t - s$. Then

$$q_s(u, u) - F_{\mathcal{L}}^1(s)^2 \langle u, u \rangle = q_t(u, u) + 2F_{\mathcal{L}}^1(s)l_t(u, u).$$

As the left side of this expression is non-negative,

$$\frac{l_t(u, u)}{q_t(u, u)} \geq -\frac{1}{2F_{\mathcal{L}}^1(s)}$$

for all $u \in \mathcal{L} \setminus \{0\}$ and the equality holds for some $u \in \mathcal{L}$. Hence $-\frac{1}{2F_{\mathcal{L}}^1(s)}$ is the smallest eigenvalue of $(Z_t^{\mathcal{L}})$, and thus necessarily equal to $\tau_1^-(t)$. In this case $s - F_{\mathcal{L}}^1(s) = t - 2F_{\mathcal{L}}^1(s) = t + \frac{1}{\tau_1^-(t)}$. Here the vector u for which equality is achieved is exactly $u = u_1^-(t)$.

Conversely, let $\tau_1^-(t)$ and $u_1^-(t)$ be as stated. Then

$$\tau_1^-(t) \leq \frac{l_t(u, u)}{q_t(u, u)}$$

for all $u \in \mathcal{L}$ with equality for $u = u_1^-(t)$. Re-arranging this expression yields

$$q_t(u, u) - \frac{1}{\tau_1^-(t)} l_t(u, u) \geq 0$$

for all $u \in \mathcal{L}$ with equality for $u = u_1^-(t)$. The substitution $t = s - \frac{1}{2\tau_1^-(t)}$ then yields

$$q_t(u, u) - \frac{1}{(2\tau_1^-(t))^2} \langle u, u \rangle \geq 0$$

for all $u \in \mathcal{L}$. The equality holds for $u = u_1^-(t)$. This expression further re-arranges as

$$\frac{|u|_s^2}{\|u\|^2} \geq \frac{1}{(2\tau_1^-(t))^2}.$$

Hence $F_{\mathcal{L}}^1(s)^2 = \frac{1}{(2\tau_1^-(t))^2}$, as needed. \square

An extension to $j \neq 1$ is now found by induction.

Theorem 6. Let $t \in \mathbb{R}$ and $1 \leq j \leq n$ be fixed.

(–) The number of negative eigenvalues $n^-(t)$ in $(Z_t^{\mathcal{L}})$ is greater than or equal to j if and only if

$$\frac{\langle Au, u \rangle}{\langle u, u \rangle} < t \quad \text{for some } u \in \mathcal{L} \ominus \text{Span}\{u_1^-(t), \dots, u_{j-1}^-(t)\}.$$

Assuming this holds true, then $\tau = \tau_j^-(t)$ and $u = u_j^-(t)$ are solutions of $(Z_t^{\mathcal{L}})$ if and only if

$$F_{\mathcal{L}}^j \left(t + \frac{1}{2\tau_j^-(t)} \right) = -\frac{1}{2\tau_j^-(t)} = \frac{|u_j^-(t)|_{t+\frac{1}{2\tau_j^-(t)}}}{\|u_j^-(t)\|}.$$

(+) The number of positive eigenvalues $n^+(t)$ in $(Z_t^{\mathcal{L}})$ is greater than or equal to j if and only if

$$\frac{\langle Au, u \rangle}{\langle u, u \rangle} > t \quad \text{for some } u \in \mathcal{L} \ominus \text{Span}\{u_1^+(t), \dots, u_{j-1}^+(t)\}.$$

Assuming this holds true, then $\tau = \tau_j^+(t)$ and $u = u_j^+(t)$ are solutions of $(Z_t^{\mathcal{L}})$ if and only if

$$F_{\mathcal{L}}^j \left(t + \frac{1}{2\tau_j^+(t)} \right) = \frac{1}{2\tau_j^+(t)} = \frac{|u_j^+(t)|_{t+\frac{1}{2\tau_j^+(t)}}}{\|u_j^+(t)\|}.$$

Proof. For $j = 1$ the statements are Lemma 5 taking into consideration (14). For $j > 1$, due to the symmetry of the eigenproblem $(Z_t^{\mathcal{L}})$, it is enough to apply again Lemma 5 by fixing $\tilde{\mathcal{L}} = \mathcal{L} \ominus \text{Span}\{u_1^{\mp}(t), \dots, u_{j-1}^{\mp}(t)\}$ as trial spaces. Note that the eigenvalues of $(Z_t^{\tilde{\mathcal{L}}})$ are those of $(Z_t^{\mathcal{L}})$ except for $\tau_1^{\mp}(t), \dots, \tau_{j-1}^{\mp}(t)$. \square

A neat procedure for finding certified spectral bounds for A , as described in [24], can now be deduced from Theorem 6. By virtue of Proposition 1 and Remark 1, this procedure turns out to be optimal in the context of the approximated counting functions discussed in Section 2, see [18, Section 6]. We summarize the core statement as follows.

Corollary 7. For all $t \in \mathbb{R}$ and $j \in \{1, \dots, n^{\pm}(t)\}$,

$$(24) \quad t + \frac{1}{\tau_j^-(t)} \leq \mathbf{n}_j^-(t) \quad \text{and} \quad \mathbf{n}_j^+(t) \leq t + \frac{1}{\tau_j^+(t)}.$$

In recent years, numerical techniques based on this statement have been designed to successfully compute eigenvalues for the radially reduced magnetohydrodynamics operator [24, 11], the Helmholtz equation [5] and the calculation of sloshing frequencies in the left definite case [4]. We will determine one such a numerical scheme for the case of the Maxwell operator in Section 6.

Remark 4. Since $\pm \frac{1}{\tau_j^{\pm}(t)} \geq \pm(\mathbf{n}_j^{\pm}(t) - t)$ in the above,

$$\hat{t}_j^- = t + \frac{1}{2\tau_j^-(t)} \leq \frac{t + \mathbf{n}_j^-(t)}{2} \leq \frac{\mathbf{n}_j^+(t) + \mathbf{n}_j^-(t)}{2} \leq \frac{\mathbf{n}_j^+(t) + t}{2} \leq t + \frac{1}{2\tau_j^+(t)} = \hat{t}_j^+.$$

Hence \hat{t}_j^{\pm} is not further from $\mathbf{n}_j^{\pm}(t)$ than it is to $\mathbf{n}_j^{\mp}(t)$. Moreover

$$\hat{t}_j^{\pm} = \frac{\mathbf{n}_j^+(t) + \mathbf{n}_j^-(t)}{2}$$

renders $t \in \sigma(A)$ and

$$\frac{1}{\tau_j^{\pm}(t)} = \mathbf{n}_j^{\pm}(t) - t.$$

This geometrical property for the solution of (23) will be relevant for our next goal, the examination of the convergence properties of the estimates (24).

4. CONVERGENCE AND ERROR ESTIMATES

Our first goal in this section will be to show that, if \mathcal{L} captures an eigenspace of A within a certain order of precision $\mathcal{O}(\varepsilon)$ as specified below, then the bounds consequence of Proposition 1 are

- a) at least within $\mathcal{O}(\varepsilon)$ from the true spectral data for any $t \in \mathbb{R}$,
- b) within $\mathcal{O}(\varepsilon^2)$ for $t \notin \sigma(A)$.

This will be the content of theorems 9 and 10, and Corollary 11. We will then show that, in turns, the estimates (24) have always residual of size $\mathcal{O}(\varepsilon^2)$ for any $t \in \mathbb{R}$. See Theorem 12. In the spectral approximation literature this property is known as optimal order of convergence/exactness, see [14, Chapter 6].

Recall Remark 2, and the assumptions 1 and 2. Below $\{\phi_j^t\}_{j=1}^m$ denotes an orthonormal set of eigenvectors of $\mathcal{E}_{[t-\mathfrak{d}_m(t), t+\mathfrak{d}_m(t)]}(A)$ which is ordered so that

$$|A - t|\phi_j^t = \mathfrak{d}_j(t)\phi_j^t \quad \text{for } j = 1, \dots, m.$$

Whenever $0 < \varepsilon_j < 1$ is small, as specified below, the trial subspace $\mathcal{L} \subset D(A)$ will be assumed to be close to $\text{Span}\{\phi_j^t\}_{j=1}^m$ in the sense that there exist $w_j^t \in \mathcal{L}$ such that

$$\begin{aligned} (A_0) \quad & \|w_j^t - \phi_j^t\| \leq \varepsilon_j \quad \text{and} \\ (A_1) \quad & |w_j^t - \phi_j^t|_t \leq \varepsilon_j. \end{aligned}$$

We have split this condition into two, in order to highlight the fact that some times only (A₁) is required. Unless otherwise specified, the index j runs from 1 to m .

From (16) it follows that the family $\{\phi_j^s\}_{j=1}^m \subset \mathcal{E}_{[t-\mathfrak{d}_m(t), t+\mathfrak{d}_m(t)]}(A)$ and the family $\{w_j^s\}_{j=1}^m \subset \mathcal{L}$ above can always be chosen piecewise constant for s in a neighbourhood of t . Moreover, they can be chosen so that jumps only occur at $s \in \sigma(A)$.

Assumption 3. *Without further mention all t -dependant vectors below will be assumed to be locally constant in t with jumps only at the spectrum of A .*

A set $\{w_j^t\}_{j=1}^m$ subject to (A₀)-(A₁) is not generally orthonormal. However, according to the next lemma, it can always be substituted by an orthonormal set, provided ε_j is small enough.

Lemma 8. *There exists a constant $C > 0$ independent of \mathcal{L} ensuring the following. If $\{w_j^t\}_{j=1}^m \subset \mathcal{L}$ is such that (A₀)-(A₁) hold for all ε_j such that*

$$\varepsilon = \sqrt{\sum_{j=1}^m \varepsilon_j^2} < \frac{1}{\sqrt{m}},$$

then there is a set $\{v_j^t\}_{j=1}^m \subset \mathcal{L}$ orthonormal in the inner product $\langle \cdot, \cdot \rangle$ such that

$$|v_j^t - \phi_j^t|_t + \|v_j^t - \phi_j^t\| < C\varepsilon.$$

Proof. As it is clear from the context, in this proof we suppress the index t on top of any vector. The desired conclusion is achieved by applying the Gram-Schmidt

procedure. Let $G = [\langle w_k, w_l \rangle]_{kl=1}^m \in \mathbb{C}^{m \times m}$ be the Gram matrix associated to $\{w_j\}$. Set

$$v_j = \sum_{k=1}^m (G^{-1/2})_{kj} w_k.$$

Then

$$\begin{aligned} \|G - I\| &\leq \sqrt{\sum_{kl=1}^m |\langle w_k, w_l \rangle - \langle \phi_k, \phi_l \rangle|^2} \\ &\leq \sqrt{2 \sum_{kl=1}^m \|w_k - \phi_k\|^2 (\|w_l\| + \|\phi_l\|)^2} \\ &\leq \sqrt{2}(2 + \varepsilon)\varepsilon. \end{aligned}$$

Since

$$\begin{aligned} \|v_j - w_j\|^2 &= \left\| \sum_{k=1}^m (G^{-1/2} - I)_{kj} w_k \right\|^2 \\ &= \sum_{kl=1}^m (G^{-1/2} - I)_{kj} \overline{(G^{-1/2} - I)_{lj}} \langle w_k, w_l \rangle \\ &= \sum_{k=1}^m (G^{-1/2} - I)_{kj} \overline{\left(\sum_{l=1}^m G_{kl} (G^{-1/2} - I)_{lj} \right)} \\ &= \sum_{k=1}^m (G^{-1/2} - I)_{kj} (G^{1/2} - G)_{jk} \\ &= \left((I - G^{1/2})^2 \right)_{jj} \end{aligned}$$

then

$$\|v_j - w_j\| \leq \|I - G^{1/2}\|.$$

As $G^{1/2}$ is a positive-definite matrix, for every $\underline{v} \in \mathbb{C}^m$ we have

$$\|(G^{1/2} + I)\underline{v}\|^2 = \|G^{1/2}\underline{v}\|^2 + 2\langle G^{1/2}\underline{v}, \underline{v} \rangle + \|\underline{v}\|^2 \geq \|\underline{v}\|^2.$$

Then $\det(I + G^{1/2}) \neq 0$ and $\|(I + G^{1/2})^{-1}\| \leq 1$. Hence

$$(25) \quad \|v_j - w_j\| \leq \|(I - G)(I + G^{1/2})^{-1}\| \leq \|I - G\| \|(I + G^{1/2})^{-1}\| \leq (2 + \varepsilon)\varepsilon.$$

Now, identify $\underline{v} = (v_1, \dots, v_m) \in \mathbb{C}^m$ with $v = \sum_{k=1}^m v_k \phi_k$. As

$$\|G^{1/2}\underline{v}\| = \left\| \sum_{j=1}^m \langle v, \phi_j \rangle w_j \right\| \geq \|v\| - \left\| \sum_{j=1}^m \langle v, \phi_j \rangle (w_j - \phi_j) \right\| \geq (1 - \varepsilon)\|v\|$$

then

$$\|G^{-1/2}\| \leq \frac{1}{1 - \varepsilon}.$$

Hence

$$\begin{aligned}
|v_j - w_j|_t &\leq \sum_{k=1}^m |(G^{-1/2} - I)_{jk}| |w_k|_t \\
&\leq \sum_{k=1}^m |(G^{-1/2} - I)_{jk}| (\varepsilon_k + \mathfrak{d}_k(t)) \\
&\leq \sum_{kl=1}^m |(G^{-1/2})_{kl}| |(G^{1/2} - I)_{lj}| (\varepsilon_k + \mathfrak{d}_k(t)) \\
(26) \quad &\leq \frac{\sqrt{m}(\varepsilon + \mathfrak{d}_m(t))(2 + \varepsilon)}{1 - \varepsilon} \varepsilon.
\end{aligned}$$

The desired conclusion follows from (25) and (26). \square

The next theorem addresses the claim *a*) made at the beginning of this section. According to Lemma 8, in order to examine the asymptotic behaviour of $F_{\mathcal{L}}^j(t)$ as $\varepsilon_j \rightarrow 0$ under the constraints (A₀)-(A₁), we can assume without loss of generality that the trial vectors w_j^t form an orthonormal set in the inner product $\langle \cdot, \cdot \rangle$.

Theorem 9. *Let $\{w_j^t\}_{j=1}^m \subset \mathcal{L}$ be a family of vectors which is orthonormal in the inner product $\langle \cdot, \cdot \rangle$ and satisfies (A₁). Then*

$$F_{\mathcal{L}}^j(t) - \mathfrak{d}_j(t) \leq \left(\sum_{k=1}^j \varepsilon_k^2 \right)^{1/2} \quad \forall j = 1, \dots, m.$$

Proof. From the min-max principle we obtain

$$\begin{aligned}
F_{\mathcal{L}}^j(t) &\leq \max_{\sum |c_k|^2=1} \left| \sum_{k=1}^j c_k w_k \right|_t \\
&\leq \max_{\sum |c_k|^2=1} \left| \sum_{k=1}^j c_k (w_k - \phi_k) \right|_t + \max_{\sum |c_k|^2=1} \left| \sum_{k=1}^j c_k \phi_k \right|_t \\
&= \max_{\sum |c_k|^2=1} \left| \sum_{k=1}^j c_k (w_k - \phi_k) \right|_t + \mathfrak{d}_j(t).
\end{aligned}$$

This gives

$$\begin{aligned}
F_{\mathcal{L}}^j(t) - \mathfrak{d}_j(t) &\leq \max_{\sum |c_k|^2=1} \sum_{k=1}^j |c_k| |w_k - \phi_k|_t \\
&\leq \max_{\sum |c_k|^2=1} \left(\sum_{k=1}^j |c_k|^2 \right)^{1/2} \left(\sum_{k=1}^j |w_k - \phi_k|_t^2 \right)^{1/2} \leq \left(\sum_{k=1}^j \varepsilon_k^2 \right)^{1/2}
\end{aligned}$$

as needed. \square

In terms of order of approximation, Theorem 9 will be superseded by Theorem 10 for $t \notin \sigma(A)$. However, if $t \in \sigma(A)$, the trial space \mathcal{L} can be chosen so that $F_{\mathcal{L}}^1(t) - \mathfrak{d}_1(t)$ is linear in ε_1 . Indeed, fixing any non-zero $u \in D(A)$ and $\mathcal{L} = \text{Span}\{u\}$, yields $F_{\mathcal{L}}^1(t) - \mathfrak{d}_1(t) = F_{\mathcal{L}}^1(t) = \varepsilon_1$. This shows that Theorem 9 is optimal, upon the presumption that t is arbitrary.

The next theorem addresses the claim *b*) made at the beginning of this section. Its proof is reminiscent of that of [23, Theorem 6.1].

Theorem 10. *Let $t \notin \sigma(A)$. Suppose that the ε_j in (A₁) are such that*

$$(27) \quad \sum_{j=1}^m \varepsilon_j^2 < \frac{\mathfrak{d}_1(t)^2}{6}.$$

Then,

$$(28) \quad F_{\mathcal{L}}^j(t) - \mathfrak{d}_j(t) \leq 3 \frac{\mathfrak{d}_j(t)}{\mathfrak{d}_1(t)^2} \sum_{k=1}^j \varepsilon_k^2 \quad \forall j = 1, \dots, m.$$

Proof. Since $t \notin \sigma(A)$, then $(D(A), q_t(\cdot, \cdot))$ is a Hilbert space. Let $P_{\mathcal{L}} : D(A) \rightarrow \mathcal{L}$ be the orthogonal projection onto \mathcal{L} with respect to the inner product $q_t(\cdot, \cdot)$, so that

$$q_t(u - P_{\mathcal{L}}u, v) = 0 \quad \forall v \in \mathcal{L}.$$

Then $|u|_t^2 = |P_{\mathcal{L}}u|_t^2 + |u - P_{\mathcal{L}}u|_t^2$ for all $u \in D(A)$ and $|u - P_{\mathcal{L}}u|_t \leq |u - v|_t$ for all $v \in \mathcal{L}$. Hence

$$(29) \quad |\phi_k - P_{\mathcal{L}}\phi_k|_t \leq \varepsilon_k \quad \forall k = 1, \dots, m.$$

Let $\mathcal{E}_j = \text{Span}\{\phi_k\}_{k=1}^j$. Define

$$\begin{aligned} \mathcal{F}_j &= \{\phi \in \mathcal{E}_j : \|\phi\| = 1\} \quad \text{and} \\ \mu_{\mathcal{L}}^j(t) &= \max_{\phi \in \mathcal{F}_j} |2 \operatorname{Re}\langle \phi, \phi - P_{\mathcal{L}}\phi \rangle - \|\phi - P_{\mathcal{L}}\phi\|^2|. \end{aligned}$$

Here $\mu_{\mathcal{L}}^j$ depends on t , as $P_{\mathcal{L}}$ does. We first show that, under hypothesis (27), $\mu_{\mathcal{L}}^j(t) < \frac{1}{2}$. Indeed, given $\phi \in \mathcal{F}_j$ we decompose it as $\phi = \sum_{k=1}^j c_k \phi_k$. Then

$$\begin{aligned} |\langle \phi, \phi - P_{\mathcal{L}}\phi \rangle| &= \left| \sum_{k=1}^j c_k \langle \phi_k, \phi - P_{\mathcal{L}}\phi \rangle \right| = \left| \sum_{k=1}^j \frac{c_k}{\mathfrak{d}_k(t)^2} q_t(\phi_k, \phi - P_{\mathcal{L}}\phi) \right| \\ &= \left| q_t \left(\sum_{k=1}^j \frac{c_k}{\mathfrak{d}_k(t)^2} \phi_k, \phi - P_{\mathcal{L}}\phi \right) \right| \\ &= \left| q_t \left(\sum_{k=1}^j \frac{c_k}{\mathfrak{d}_k(t)^2} (\phi_k - P_{\mathcal{L}}\phi_k), \phi - P_{\mathcal{L}}\phi \right) \right| \\ (30) \quad &\leq \left| \sum_{k=1}^j \frac{c_k}{\mathfrak{d}_k(t)^2} (\phi_k - P_{\mathcal{L}}\phi_k) \right|_t \left| \sum_{k=1}^j c_k (\phi_k - P_{\mathcal{L}}\phi_k) \right|_t. \end{aligned}$$

For each multiplying term in the latter expression, the triangle and Cauchy-Schwarz's inequalities yield (take $\alpha_k = c_k$ or $\alpha_k = \frac{c_k}{\mathfrak{d}_k(t)^2}$)

$$\begin{aligned} \left| \sum_{k=1}^j \alpha_k (\phi_k - P_{\mathcal{L}}\phi_k) \right|_t &\leq \sum_{k=1}^j |\alpha_k| |\phi_k - P_{\mathcal{L}}\phi_k|_t \\ (31) \quad &\leq \left(\sum_{k=1}^j |\alpha_k|^2 \right)^{1/2} \left(\sum_{k=1}^j |\phi_k - P_{\mathcal{L}}\phi_k|_t^2 \right)^{1/2}. \end{aligned}$$

Then

$$\begin{aligned}
(32) \quad |2 \operatorname{Re}\langle \phi, \phi - P_{\mathcal{L}}\phi \rangle| &\leq 2 \left(\sum_{k=1}^j \frac{|c_k|^2}{\mathfrak{d}_k(t)^4} \right)^{1/2} \left(\sum_{k=1}^j |c_k|^2 \right)^{1/2} \sum_{k=1}^j \varepsilon_k^2 \\
&\leq \frac{2}{\mathfrak{d}_1(t)^2} \sum_{k=1}^j \varepsilon_k^2
\end{aligned}$$

for all $\phi \in \mathcal{F}_j$.

The other term in the expression for $\mu_{\mathcal{L}}^j(t)$ has an upper bound found as follows. According to the min-max principle

$$(33) \quad \|\phi - P_{\mathcal{L}}\phi\|^2 \leq \frac{1}{\mathfrak{d}_1(t)^2} q_t(\phi - P_{\mathcal{L}}\phi, \phi - P_{\mathcal{L}}\phi).$$

Therefore, by repeating analogous steps as in (30) and (31), we get

$$\begin{aligned}
(34) \quad \|\phi - P_{\mathcal{L}}\phi\|^2 &\leq \frac{1}{\mathfrak{d}_1(t)^2} \sum_{k=1}^j c_k q_t(\phi_k - P_{\mathcal{L}}\phi_k, \phi - P_{\mathcal{L}}\phi) \\
&= q_t \left(\sum_{k=1}^j \frac{c_k}{\mathfrak{d}_1(t)^2} (\phi_k - P_{\mathcal{L}}\phi_k), \phi - P_{\mathcal{L}}\phi \right) \\
&= q_t \left(\sum_{k=1}^j \frac{c_k}{\mathfrak{d}_1(t)^2} (\phi_k - P_{\mathcal{L}}\phi_k), \sum_{l=1}^j c_l (\phi_l - P_{\mathcal{L}}\phi_l) \right) \\
&\leq \frac{1}{\mathfrak{d}_1(t)^2} \sum_{k=1}^j \varepsilon_k^2.
\end{aligned}$$

Hence, from (32) and (34),

$$(35) \quad \mu_{\mathcal{L}}^j(t) \leq \frac{3}{\mathfrak{d}_1(t)^2} \sum_{k=1}^j \varepsilon_k^2 < \frac{1}{2}$$

as a consequence of (27).

Next, observe that $\dim(P_{\mathcal{L}}\mathcal{E}_j) = j$. Indeed $P_{\mathcal{L}}\psi = 0$ for $\|\psi\| = 1$ would imply

$$\mu_{\mathcal{L}}^j(t) \geq |2 \operatorname{Re}\langle \psi, \psi - P_{\mathcal{L}}\psi \rangle - \|\psi - P_{\mathcal{L}}\psi\|^2| = \|\psi\|^2 = 1,$$

which would contradict the fact that $\mu_{\mathcal{L}}^j(t) < 1$. Then,

$$F_{\mathcal{L}}^j(t)^2 \leq \max_{u \in P_{\mathcal{L}}\mathcal{E}_j} \frac{|u|_t^2}{\|u\|^2} = \max_{\phi \in \mathcal{E}_j} \frac{|P_{\mathcal{L}}\phi|_t^2}{\|P_{\mathcal{L}}\phi\|^2} = \max_{\phi \in \mathcal{F}_j} \frac{|P_{\mathcal{L}}\phi|_t^2}{\|P_{\mathcal{L}}\phi\|^2}.$$

As

$$\|P_{\mathcal{L}}\phi\|^2 = \|\phi\|^2 - 2 \operatorname{Re}\langle \phi, \phi - P_{\mathcal{L}}\phi \rangle + \|\phi - P_{\mathcal{L}}\phi\|^2 \geq 1 - \mu_{\mathcal{L}}^j(t),$$

we get

$$(36) \quad F_{\mathcal{L}}^j(t)^2 \leq \max_{\phi \in \mathcal{F}_j} \frac{|\phi|_t^2}{1 - \mu_{\mathcal{L}}^j(t)} = \max_{\sum |c_k|^2 = 1} \frac{\sum_{k=1}^j |c_k|^2 \mathfrak{d}_k(t)^2}{1 - \mu_{\mathcal{L}}^j(t)} = \frac{\mathfrak{d}_j(t)^2}{1 - \mu_{\mathcal{L}}^j(t)}.$$

Finally, (36) and (35) yield

$$\begin{aligned}
F_{\mathcal{L}}^j(t)^2 - \mathfrak{d}_j(t)^2 &\leq \frac{\mu_{\mathcal{L}}^j(t)}{1 - \mu_{\mathcal{L}}^j(t)} \mathfrak{d}_j(t)^2 \\
&\leq 2\mu_{\mathcal{L}}^j(t) \mathfrak{d}_j(t)^2 \\
(37) \qquad \qquad \qquad &\leq 2 \frac{3}{\mathfrak{d}_1(t)^2} \mathfrak{d}_j(t)^2 \sum_{k=1}^j \varepsilon_k^2.
\end{aligned}$$

The proof is completed by observing that $F_{\mathcal{L}}^j(t) + \mathfrak{d}_j(t) \geq 2\mathfrak{d}_j(t)$. \square

As the next corollary shows, a quadratic order of decrease for $F_{\mathcal{L}}^j(t) - \mathfrak{d}_j(t)$ is prevented for $t \in \sigma(A)$ in the context of theorems 9 and 10, only for j up to $\dim \mathcal{E}_t(A)$.

Corollary 11. *Let $t \in \sigma_{\text{disc}}(A)$, $\ell = 1 + \dim \mathcal{E}_t(A)$ and $k \in \{\ell, \dots, m\}$. Let*

$$\alpha_k(t) = \frac{1}{4} \min \{ |\mathfrak{d}_l(t) - \mathfrak{d}_{l-1}(t)| : \mathfrak{d}_l(t) \neq \mathfrak{d}_{l-1}(t), l = \ell, \dots, k \} > 0.$$

There exists $\varepsilon > 0$ independent of k ensuring the following. If (A_1) holds true for $\sqrt{\sum_{j=1}^m \varepsilon_j^2} < \varepsilon$, then

$$F_{\mathcal{L}}^k(t) - \mathfrak{d}_k(t) \leq 3 \frac{\mathfrak{d}_k(t)}{\alpha_k(t)^2} \sum_{j=1}^k \varepsilon_j^2.$$

Proof. Without loss of generality we assume that $t + \mathfrak{d}_k(t) \in \sigma(A)$. Otherwise $t - \mathfrak{d}_k(t) \in \sigma(A)$ and the proof is analogous to the one presented below.

Let $\tilde{t} = t + \alpha_k(t)$. Then $\tilde{t} \notin \sigma(A)$ and $t + \mathfrak{d}_k(t) = \tilde{t} + \mathfrak{d}_k(\tilde{t})$. Since the map $s \mapsto s + F_{\mathcal{L}}^j(s)$ is non-decreasing as a consequence of Proposition 1, Theorem 10 applied at \tilde{t} yields

$$\begin{aligned}
F_{\mathcal{L}}^k(t) - \mathfrak{d}_k(t) &= t + F_{\mathcal{L}}^k(t) - (t + \mathfrak{d}_k(t)) \leq \tilde{t} + F_{\mathcal{L}}^k(\tilde{t}) - (\tilde{t} + \mathfrak{d}_k(\tilde{t})) \\
&= F_{\mathcal{L}}^k(\tilde{t}) - \mathfrak{d}_k(\tilde{t}) \leq 3 \frac{\mathfrak{d}_k(\tilde{t})}{\mathfrak{d}_1(\tilde{t})^2} \sum_{j=1}^k \varepsilon_j^2 \leq 3 \frac{\mathfrak{d}_k(t)}{\alpha_k(t)^2} \sum_{j=1}^k \varepsilon_j^2
\end{aligned}$$

as needed. \square

For the final part of this section, we are now able to formulate a precise statements on the convergence of the method of Zimmermann and Mertins. Theorem 12 below improves upon two crucial aspects of a similar result established in [11, Lemma 2]. It allows $j > 1$ and it allows $t \in \sigma(A)$. These two improvements are essential in order to obtain sharp bounds for those eigenvalues of the Maxwell operator which are either degenerate or form a tight cluster.

Remark 5. The constants $\tilde{\varepsilon}_t$ and C_t^{\pm} below do have a dependance on t that may be determined explicitly from Theorem 10, Corollary 11 and the proof of Theorem 12. Despite of the fact that they can deteriorate as t approaches the isolated eigenvalues of A and they can have jumps precisely at these points, they may be chosen locally independent of t in compacts outside the spectrum. This has an impact on practical implementations of the computational method to be described in Section 6 which we do not fully understand at present. Our numerical tests in Section 7 indicate

that the best results are achieved by choosing t relatively far from the spectral point being approximated.

Set

$$\begin{aligned}\nu_j^-(t) &= \sup\{s < t : \text{Tr } \mathbf{1}_{(s,t)}(A) \geq j\} \\ \nu_j^+(t) &= \inf\{s > t : \text{Tr } \mathbf{1}_{(t,s)}(A) \geq j\}.\end{aligned}$$

Note that these are the spectral points of A which are strictly to the left and strictly to the right of t respectively. The inequality $\nu_j^\pm(t) \neq n_j^\pm(t)$ only occurs when t is an eigenvalue. In view of (15), $\delta^\pm(t) = |t - \nu_1^\pm(t)|$.

Theorem 12. *Let $J \subset \mathbb{R}$ be a bounded open segment such that $J \cap \sigma(A) \subseteq \sigma_{\text{disc}}(A)$. Let $\{\phi_k\}_{k=1}^{\tilde{m}}$ be a family of eigenvectors of A such that $\text{Span}\{\phi_k\}_{k=1}^{\tilde{m}} = \mathcal{E}_J(A)$. For fixed $t \in J$, there exist constants $\tilde{\varepsilon}_t > 0$ and $C_t^\pm > 0$ independent of the trial space \mathcal{L} , ensuring the following. If there are $\{w_j\}_{j=1}^{\tilde{m}} \subset \mathcal{L}$ such that*

$$(38) \quad \left(\sum_{j=1}^{\tilde{m}} \|w_j - \phi_j\|^2 + |w_j - \phi_j|_t^2 \right)^{1/2} \leq \varepsilon < \tilde{\varepsilon}_t,$$

then

$$\left| \nu_j^\pm(t) - \left(t + \frac{1}{\tau_j^\pm(t)} \right) \right| \leq C_t^\pm \varepsilon^2$$

for all $j \leq n^\pm(t)$ such that $\nu_j^\pm(t) \in J$.

Proof. We focus on the case of the plus sign, as the one with the minus sign is completely analogous. The hypotheses ensure that the number of indices $j \leq n^\pm(t)$ such that $\nu_j^\pm(t) \in J$ never exceeds \tilde{m} . Therefore this condition in the conclusion of the theorem is consistent.

Let

$$m(t) = \max\{m \in \mathbb{N} : [t - \mathfrak{d}_m(t), t + \mathfrak{d}_m(t)] \subset J\}.$$

Recall the Assumption 1 and the Remark 3. The hypothesis on \mathcal{L} guarantees that (A₀)-(A₁) hold true for $m = m(t)$ and $\left(\sum_{j=1}^{m(t)} \varepsilon_j^2 \right)^{1/2} < \varepsilon$. By combining Lemma 8, Theorem 9 and the fact that we can pick $\{w_j^t\}_{j=1}^{m(t)} \subseteq \{w_k\}_{k=1}^{\tilde{m}}$, there exists $\tilde{\varepsilon}_t > 0$ small enough, such that (38) yields

$$(39) \quad F_{\mathcal{L}}^j(s) - \mathfrak{d}_j(s) \leq \frac{\nu_1^+(t) - t}{2} \quad \forall j = 1, \dots, \tilde{m} \quad \text{and} \quad s \in J.$$

Let j be such that $\nu_j^\pm(t) \in J$. Since $t + \alpha - \nu_1^+(t) \leq \nu_j^+(t) - (\alpha + t)$ for all $0 \leq \alpha \leq \frac{\nu_j^+(t) + \nu_1^+(t)}{2} - t$, then

$$\mathfrak{d}_j(s) = \nu_j^+(t) - s \quad \forall s \in \left[\frac{t + \nu_j^+(t)}{2}, \frac{\nu_1^+(t) + \nu_j^+(t)}{2} \right].$$

Let

$$g(\alpha) = \alpha - F_{\mathcal{L}}^j(t + \alpha).$$

Then g is an increasing function of α and $g(0) = -F_{\mathcal{L}}^j(t) < 0$. For the strict inequality in the latter, recall Assumption 2. Moreover, according to (39)

$$\begin{aligned} g\left(\frac{\nu_j^+(t) + \nu_1^+(t)}{2} - t\right) &= \nu_1^+(t) - t + \frac{\nu_j^+(t) - \nu_1^+(t)}{2} - F_{\mathcal{L}}^j\left(\frac{\nu_j^+(t) + \nu_1^+(t)}{2}\right) \\ &= \nu_1^+(t) - t + \mathfrak{d}_j\left(\frac{\nu_j^+(t) + \nu_1^+(t)}{2}\right) - F_{\mathcal{L}}^j\left(\frac{\nu_j^+(t) + \nu_1^+(t)}{2}\right) \\ &\geq \nu_1^+(t) - t - \frac{\nu_1^+(t) - t}{2} > 0. \end{aligned}$$

Hence, the Mean Value Theorem ensures the existence of $\tilde{\alpha} \in \left(0, \frac{\nu_1^+(t) + \nu_j^+(t)}{2} - t\right)$ such that $\tilde{\alpha} = F_{\mathcal{L}}^j(t + \tilde{\alpha})$. According to Theorem 6 (+), $\tilde{\alpha}$ is unique and $\tilde{\alpha} = \frac{1}{2\tau_j^+(t)}$.

The proof is now completed as follows. By virtue of Remark 4,

$$\hat{t}_j^+(t) = t + \frac{1}{2\tau_j^+(t)} \in \left(\frac{t + \nu_j^+(t)}{2}, \frac{\nu_1^+(t) + \nu_j^+(t)}{2}\right) \quad \text{and} \quad F_{\mathcal{L}}^j(\hat{t}_j^+(t)) = \frac{1}{2\tau_j^+(t)}.$$

Then, Theorem 10 or Corollary 11, as appropriate, ensure the existence of $C_t^+ > 0$ yielding

$$\nu_j^+(t) - \left(t + \frac{1}{\tau_j^+(t)}\right) = F_{\mathcal{L}}^j(\hat{t}_j^+) - \mathfrak{d}_j(\hat{t}_j^+) \leq C_t^+ \sum_{k=1}^j \varepsilon_k^2 < \varepsilon^2,$$

as needed. \square

We conclude this section with a result on convergence of eigenfunctions.

Corollary 13. *Let $J \subset \mathbb{R}$ be a bounded open segment such that $J \cap \sigma(A) \subseteq \sigma_{\text{disc}}(A)$. Let $\{\phi_k\}_{k=1}^m$ be a family of eigenvectors of A such that $\text{Span}\{\phi_k\}_{k=1}^m = \mathcal{E}_J(A)$. For fixed $t \in J$, there exist constants $\tilde{\varepsilon}_t > 0$ and $C_t^\pm > 0$ independent of the trial space \mathcal{L} , ensuring the following. If there are $\{w_j\}_{j=1}^m \subset \mathcal{L}$ guaranteeing the validity of (38), then for all $j \leq n^\pm(t)$ such that $\nu_j^\pm(t) \in J$ there exist $\psi_j^{\varepsilon^\pm} \in \mathcal{E}_{\{\nu_j^-(t), \nu_j^+(t)\}}(A)$ such that*

$$|u_j^\pm(t) - \psi_j^{\varepsilon^\pm}|_t + \|u_j^\pm(t) - \psi_j^{\varepsilon^\pm}\| \leq C_t^\pm \varepsilon.$$

Proof. Fix $t \in J$. By virtue of Theorem 6, $u_j^\pm(t) = u_j^{\hat{t}_j^\pm}$ in the notation for eigenvectors employed in Proposition 3. The claimed conclusion is a consequence of the latter combined with Theorem 10 or Corollary 11, as appropriate. \square

5. THE FINITE ELEMENT METHOD FOR THE MAXWELL EIGENVALUE PROBLEM

Let Ω be an open subset of \mathbb{R}^3 . Below $\mathcal{D}(\Omega)$ denotes the infinitely differentiable test functions with compact support in Ω . The inner product of $L^2(\Omega)$ is $\langle \cdot, \cdot \rangle_\Omega$ and its norm $\|\cdot\|_{0,\Omega}$. The Sobolev space of order m is $\mathcal{H}^m(\Omega)$ and its norm is $\|\cdot\|_{m,\Omega}$. We do not distinguish in the notation between products and norms of scalar functions or vector fields with components in these linear spaces.

We define rigorously the domain of the operator \mathcal{M} associated to the eigenvalue problem (1) by following closely the ideas of the work [7]. Let

$$\mathcal{H}(\text{curl}; \Omega) = \{\mathbf{u} \in [L^2(\Omega)]^3 : \text{curl } \mathbf{u} \in [L^2(\Omega)]^3\}$$

equipped with the norm

$$(40) \quad \|\mathbf{u}\|_{\text{curl},\Omega}^2 = \|\mathbf{u}\|_{0,\Omega}^2 + \|\text{curl } \mathbf{u}\|_{0,\Omega}^2.$$

Let \mathcal{R}_{\max} denote the operator defined by the expression ‘‘curl’’ acting on the domain $D(\mathcal{R}_{\max}) = \mathcal{H}(\text{curl}; \Omega)$, the maximal domain. Let

$$\mathcal{R}_{\min} = \mathcal{R}_{\max}^* = \overline{\mathcal{R}_{\max} \upharpoonright [\mathcal{D}(\Omega)]^3}.$$

The domain of \mathcal{R}_{\min} is

$$\begin{aligned} D(\mathcal{R}_{\min}) &= \mathcal{H}_0(\text{curl}; \Omega) \\ &= \{\mathbf{u} \in \mathcal{H}(\text{curl}; \Omega) : \langle \text{curl } \mathbf{u}, \mathbf{v} \rangle_{\Omega} = \langle \mathbf{u}, \text{curl } \mathbf{v} \rangle_{\Omega} \quad \forall \mathbf{v} \in \mathcal{H}(\text{curl}; \Omega)\}. \end{aligned}$$

By virtue of Green’s identity for the rotational (see e.g. [20, Theorem I.2.11]), if Ω is Lipschitz in the sense of [1, Notation 2.1], then $\mathbf{u} \in \mathcal{H}_0(\text{curl}; \Omega)$ if and only if $\mathbf{u} \in \mathcal{H}(\text{curl}; \Omega)$ and $\mathbf{u} \times \mathbf{n} = \mathbf{0}$ on $\partial\Omega$.

Let

$$\mathcal{M}_1 = \begin{pmatrix} 0 & i\mathcal{R}_{\max} \\ -i\mathcal{R}_{\min} & 0 \end{pmatrix}$$

on the domain

$$(41) \quad D(\mathcal{M}_1) = D(\mathcal{R}_{\min}) \times D(\mathcal{R}_{\max}) \subset [L^2(\Omega)]^6.$$

As \mathcal{R}_{\max} and \mathcal{R}_{\min} are mutually adjoints, $\mathcal{M}_1 : D(\mathcal{M}_1) \rightarrow [L^2(\Omega)]^6$ is a self-adjoint operator, [7, Lemma 1.2]. Now, write the system (1) as

$$\begin{pmatrix} \epsilon^{-1/2} & 0 \\ 0 & \mu^{-1/2} \end{pmatrix} \begin{pmatrix} 0 & i \text{curl} \\ -i \text{curl} & 0 \end{pmatrix} \begin{pmatrix} \epsilon^{-1/2} & 0 \\ 0 & \mu^{-1/2} \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{E}} \\ \tilde{\mathbf{H}} \end{pmatrix} = \omega \begin{pmatrix} \tilde{\mathbf{E}} \\ \tilde{\mathbf{H}} \end{pmatrix}$$

with unknowns $(\tilde{\mathbf{E}}, \tilde{\mathbf{H}}) = (\epsilon^{1/2} \mathbf{E}, \mu^{1/2} \mathbf{H})$. Let

$$\mathcal{P} = \text{diag}[\epsilon^{1/2} I_{3 \times 3}, \mu^{1/2} I_{3 \times 3}]$$

be the self-adjoint operator acting on $[L^2(\Omega)]^6$ given by the permittivity and permeability. The constraint (2) ensures that \mathcal{P} is bounded and invertible with

$$\mathcal{P}^{-1} = \text{diag}[\epsilon^{-1/2} I_{3 \times 3}, \mu^{-1/2} I_{3 \times 3}].$$

Define $\mathcal{M} = \mathcal{P}^{-1} \mathcal{M}_1 \mathcal{P}^{-1}$ on the dense domain $D(\mathcal{M}) = \mathcal{P}(D(\mathcal{M}_1))$. Then \mathcal{M} is a self-adjoint operator and its eigenvalues correspond exactly with the angular frequencies in (1). Every eigenfunction $(\tilde{\mathbf{E}}, \tilde{\mathbf{H}})^t \neq 0$ of \mathcal{M} will produce a corresponding field phasor $(\mathbf{E}, \mathbf{H})^t = \mathcal{P}^{-1}(\tilde{\mathbf{E}}, \tilde{\mathbf{H}})^t \neq 0$ satisfying (1) and vice-versa.

Assumption 4. *Here and everywhere below we assume that the non-zero spectrum of \mathcal{M}_1 is purely discrete and it does not accumulate at $\omega = 0$. This hypothesis can be verified whenever Ω is a polyhedron with Lipschitz boundary for example, see [22, Corollary 3.49] and [7, Lemma 1.3]. A more systematic analysis of the properties of \mathcal{M} on more general regions Ω will be carried out elsewhere [3].*

Suppose that Ω is a polyhedron. We may consider applying the framework of Section 3 for $A = \mathcal{M}$ as follows. Fix $\{\mathcal{T}_h\}_{h>0}$ a family of shape-regular triangulations of $\bar{\Omega}$ [19], where the elements $K \in \mathcal{T}_h$ are simplexes with diameter h_K such that $h = \max_{K \in \mathcal{T}_h} h_K$. For $r \geq 1$, let

$$\begin{aligned} \mathbf{V}_h^r &= \{\mathbf{v}_h \in [C^0(\bar{\Omega})]^3 : \mathbf{v}_h|_K \in [\mathbb{P}_r(K)]^3 \quad \forall K \in \mathcal{T}_h\} \\ \mathbf{V}_{h,0}^r &= \{\mathbf{v}_h \in \mathbf{V}_h^r : \mathbf{v}_h \times \mathbf{n} = \mathbf{0} \text{ on } \partial\Omega\}. \end{aligned}$$

Then

$$(42) \quad \mathcal{L}_h = \mathbf{V}_{h,0}^r \times \mathbf{V}_h^r \subset \mathbf{D}(\mathcal{M}_1)$$

and

$$(43) \quad \tilde{\mathcal{L}}_h = \mathcal{P}\mathcal{L}_h \subset \mathbf{D}(\mathcal{M})$$

are finite element spaces of isotropic and anisotropic media, respectively. Recall that \mathcal{P} is bounded and invertible as a consequence of (2).

By virtue of [22, Theorem 3.26] and the fact that $\mathcal{H}_0(\text{curl}; \Omega)$ is the closure in the curl norm of $C_0^\infty(\Omega)$, the family \mathcal{L}_h is dense in $\mathbf{D}(\mathcal{M}_1)$. That is, for any $(\mathbf{F}, \mathbf{G})^t \in \mathbf{D}(\mathcal{M}_1)$ there exists a sequence $\{(\mathbf{F}_h, \mathbf{G}_h)^t\}_{h>0}$ such that $(\mathbf{F}_h, \mathbf{G}_h)^t \in \mathcal{L}_h$ and

$$(44) \quad \lim_{h \rightarrow 0} \left(\|\mathbf{F} - \mathbf{F}_h\|_{\text{curl}, \Omega} + \|\mathbf{G} - \mathbf{G}_h\|_{\text{curl}, \Omega} \right) = 0.$$

In turns, this implies that for all $(\tilde{\mathbf{F}}, \tilde{\mathbf{G}})^t = \mathcal{P}(\mathbf{F}, \mathbf{G})^t \in \mathbf{D}(\mathcal{M})$, there exists a family $\{(\tilde{\mathbf{F}}_h, \tilde{\mathbf{G}}_h)^t\}_{h>0} \subset \tilde{\mathcal{L}}_h$ such that

$$(45) \quad \lim_{h \rightarrow 0} \left(\left\| \mathcal{M} \begin{pmatrix} \tilde{\mathbf{F}} - \tilde{\mathbf{F}}_h \\ \tilde{\mathbf{G}} - \tilde{\mathbf{G}}_h \end{pmatrix} \right\|_{0, \Omega} + \left\| \begin{pmatrix} \tilde{\mathbf{F}} - \tilde{\mathbf{F}}_h \\ \tilde{\mathbf{G}} - \tilde{\mathbf{G}}_h \end{pmatrix} \right\|_{0, \Omega} \right) = 0.$$

Let \mathcal{I}_h denote the Lagrange interpolator on \mathcal{L}_h , [19]. Under the condition of regularity $(\mathbf{F}, \mathbf{G})^t \in \mathcal{H}^{r+1}(\Omega)^6$,

$$(46) \quad \|\mathbf{F} - \mathcal{I}_h(\mathbf{F})\|_{\text{curl}, \Omega} + \|\mathbf{G} - \mathcal{I}_h(\mathbf{G})\|_{\text{curl}, \Omega} \leq C_r h^r (\|\mathbf{F}\|_{r+1, \Omega} + \|\mathbf{G}\|_{r+1, \Omega})$$

for a suitable constant $C_r > 0$. Hence, there also exists a constant $\tilde{C}_r(\epsilon, \mu) > 0$, such that

$$(47) \quad \left\| \mathcal{M} \begin{pmatrix} \tilde{\mathbf{F}} - \mathcal{I}_h(\tilde{\mathbf{F}}) \\ \tilde{\mathbf{G}} - \mathcal{I}_h(\tilde{\mathbf{G}}) \end{pmatrix} \right\|_{0, \Omega} + \left\| \begin{pmatrix} \tilde{\mathbf{F}} - \tilde{\mathbf{F}}_h \\ \tilde{\mathbf{G}} - \tilde{\mathbf{G}}_h \end{pmatrix} \right\|_{0, \Omega} \leq \tilde{C}_r(\epsilon, \mu) h^r.$$

As a consequence of Theorem 12 and Corollary 13, the estimates (45) and (47) lead to precise convergence and error estimates for the method of Section 3 in the case $A = \mathcal{M}$ and $\mathcal{L} = \tilde{\mathcal{L}}_h$. We summarize the corresponding statements in two main theorems.

Theorem 14. *Let $J \subset \mathbb{R}$ be a bounded open segment such that $0 \notin J$. Let $t \in J$. Let $\tau_{j,h}^+(t)$ and $\tau_{j,h}^-(t)$ be the corresponding positive and negative eigenvalues of $(\mathbf{Z}_t^\mathcal{L})$ for $\mathcal{L} = \tilde{\mathcal{L}}_h$. Then, for every j such that $\nu_j^\pm(t) \in J$,*

$$\lim_{h \rightarrow 0} \left| \left(t + \frac{1}{\tau_{j,h}^\pm(t)} \right) - \nu_j^\pm(t) \right| = 0.$$

Moreover, if in addition $\mathcal{P}^{-1}\mathcal{E}_J(\mathcal{M}) \subseteq \mathcal{H}^{r+1}(\Omega)^6$, then there exist $C_t^\pm > 0$ such that

$$(48) \quad \left| \left(t + \frac{1}{\tau_{j,h}^\pm(t)} \right) - \nu_j^\pm(t) \right| \leq C_t^\pm h^{2r}$$

for h sufficiently small and j such that $\nu_j^\pm(t) \in J$.

For $(\tilde{\mathbf{F}}, \tilde{\mathbf{G}})^t \in \mathbf{D}(\mathcal{M})$ and a subspace $\mathcal{E} \subseteq \mathbf{D}(\mathcal{M})$, let

$$\text{dist}_{\mathcal{M}}[(\tilde{\mathbf{F}}, \tilde{\mathbf{G}}), \mathcal{E}] = \inf_{(\mathbf{X}, \mathbf{Y})^t \in \mathcal{E}} \left[\left\| \mathcal{M} \begin{pmatrix} \tilde{\mathbf{F}} - \mathbf{X} \\ \tilde{\mathbf{G}} - \mathbf{Y} \end{pmatrix} \right\|_{0,\Omega} + \left\| \begin{pmatrix} \tilde{\mathbf{F}} - \mathbf{X} \\ \tilde{\mathbf{G}} - \mathbf{Y} \end{pmatrix} \right\|_{0,\Omega} \right].$$

Theorem 15. *Assume the same hypotheses as in Theorem 14. Let*

$$(\tilde{\mathbf{E}}_{j,h}^{\pm}(t), \tilde{\mathbf{H}}_{j,h}^{\pm}(t))^t \in \tilde{\mathcal{L}}_h$$

be the corresponding normalized eigenvectors of the eigenvalue problem $(\mathbf{Z}_t^{\tilde{\mathcal{L}}_h})$. Then, for every j such that $\nu_j^{\pm}(t) \in J$,

$$\lim_{h \rightarrow 0} \text{dist}_{\mathcal{M}}[(\tilde{\mathbf{E}}_{j,h}^{\pm}(t), \tilde{\mathbf{H}}_{j,h}^{\pm}(t)), \mathcal{E}_{\{\nu_j^-(t), \nu_j^+(t)\}}(\mathcal{M})] = 0.$$

Moreover, if in addition $\mathcal{P}^{-1}\mathcal{E}_J(\mathcal{M}) \subseteq \mathcal{H}^{r+1}(\Omega)^6$, then there exist $C_t^{\pm} > 0$ such that

$$\text{dist}_{\mathcal{M}}[(\tilde{\mathbf{E}}_{j,h}^{\pm}(t), \tilde{\mathbf{H}}_{j,h}^{\pm}(t)), \mathcal{E}_{\{\nu_j^-(t), \nu_j^+(t)\}}(\mathcal{M})] \leq C_t^{\pm} h^r$$

for h sufficiently small and j such that $\nu_j^{\pm}(t) \in J$.

Theorems 14 and 15 above have various consequences for the numerical calculation of the eigenfrequencies associated to the resonant cavity problem which are worth highlighting. Note that convergence and absence of spectral pollution are guaranteed, despite of the fact that \mathcal{L}_h is a spaces of nodal finite elements with no particular mesh structure. These convergence properties are constrained to extremely mild assumptions on the coefficients ϵ and μ . Moreover, the order of approximation achieved is optimal in the context of the finite elements chosen.

Our analysis above relies on the regularity of the eigenspaces associated to the interval J only. This opens the possibility of approximating eigenvalues associated to regular eigenfunctions with high accuracy, if a priori information about their location is at hand. Refer to the numerical results below for concrete examples on this matter.

The discussion above was restricted finite elements of Lagrange type with the sole purpose of illustrating a concrete implementation. Analogous approximation results hold true for other choices of trial subspaces (made out of standard finite elements or otherwise) as long as they form a dense family in $\mathbf{D}(\mathcal{M})$. A control in the order of convergence will be achieved in a similar fashion, as long as interpolation estimates are available.

6. THE NUMERICAL STRATEGY IN A NUTSHELL

We now describe a certified numerical scheme for computing the eigenvalues of \mathcal{M} which is based on Corollary 7. In an asymptotic regime, as specified below, this scheme provides small intervals which are guaranteed to contain spectral points. Its convergence will be deduced from Theorem 14.

Let $t > 0$. Let $\mathcal{L} = \tilde{\mathcal{L}}_h$ as in (42)-(43) satisfy (14). Bounds for the eigenvalues of \mathcal{M} in a vicinity of t , can be found from (24). The inverse residuals $\tau_j^{\mp}(t)$ in (24) can be computed by solving $(\mathbf{Z}_t^{\tilde{\mathcal{L}}_h})$ as follows. Let $\{b_1, \dots, b_{n(h)}\}$ be a basis of \mathcal{L}_h . Let $B_t, K_t \in \mathbb{C}^{n(h) \times n(h)}$ be determined by

$$\begin{aligned} [B_t]_{jk} &= \langle (\mathcal{P}^{-1}\mathcal{M}_1 - t\mathcal{P})b_j, (\mathcal{P}^{-1}\mathcal{M}_1 - t\mathcal{P})b_k \rangle \quad \text{and} \\ [K_t]_{jk} &= \langle (\mathcal{P}^{-1}\mathcal{M}_1 - t\mathcal{P})b_j, \mathcal{P}b_k \rangle. \end{aligned}$$

Then $\tau_j^\mp(t) = \eta_\mp^{-1}$ where η_\mp is the negative(-)/positive(+) eigenvalue of the pencil $B_t - \eta K_t$ which is in the j th place among those closer to 0.

Denote by $0 < t_{\text{up}} < t_{\text{low}}$ the corresponding position t set for computing upper and lower bounds by means of $\tau_j^-(t_{\text{low}})$ and $\tau_j^+(t_{\text{up}})$, respectively. Since \mathcal{M} is strongly indefinite and $\tilde{\mathcal{L}}_h$ are dense in the graph norm of $D(\mathcal{M})$ for suitable sub-families of mesh, we can always assume that the trial spaces are chosen such that

$$(49) \quad \min_{u \in \tilde{\mathcal{L}}_h} \frac{\langle \mathcal{M}u, u \rangle}{\langle u, u \rangle} < t_{\text{up}} \quad \text{and} \quad t_{\text{low}} < \max_{u \in \tilde{\mathcal{L}}_h} \frac{\langle \mathcal{M}u, u \rangle}{\langle u, u \rangle}.$$

Recall the condition (14).

The following procedure aims at computing intervals of enclosure for the eigenvalues of \mathcal{M} which lie in the segment $(t_{\text{up}}, t_{\text{low}})$ for a prescribed tolerance set by the parameter $\delta > 0$. According to Lemma 16 below, these intervals will be certified in the regime $\delta \rightarrow 0$.

Procedure 1.

Input.

- Initial $t_{\text{up}} > 0$.
- Initial $t_{\text{low}} > t_{\text{up}}$ such that $t_{\text{low}} - t_{\text{up}}$ is fairly large.
- A sub-family \mathcal{F} of finite element spaces $\tilde{\mathcal{L}}_h$ as in (42)-(43), dense in the graph norm of $D(\mathcal{M})$ as $h \rightarrow 0$.
- A tolerance $\delta > 0$ fairly small compared with $t_{\text{low}} - t_{\text{up}}$.

Output.

- A prediction $\tilde{m}(\delta) \in \mathbb{N}$ of $\text{Tr} \mathbf{1}_{(t_{\text{up}}, t_{\text{low}})}(\mathcal{M})$.
- Predictions $\omega_j^\pm(\delta)$ of the endpoints of enclosures for the eigenvalues in $\sigma(\mathcal{M}) \cap (t_{\text{up}}, t_{\text{low}})$, such that $0 < \omega_j^+(\delta) - \omega_j^-(\delta) < \delta$ for $j = 1, \dots, \tilde{m}(\delta)$.

Steps.

- a) Set initial $\mathcal{L} = \tilde{\mathcal{L}}_h \in \mathcal{F}$.
- b) While

$$\omega_{j,h}^+ - \omega_{j,h}^- \geq \delta \text{ or } \omega_{j,h}^- > \omega_{j,h}^+ \text{ for some } j = 1, \dots, \tilde{m},$$

do c) - e).

- c) Compute

$$\omega_{j,h}^+ = t_{\text{up}} + \frac{1}{\tau_{j^+}^+(t_{\text{up}})} \quad \text{for } j = 1, \dots, \tilde{m}_{\text{up}}$$

where \tilde{m}_{up} is such that all $\omega_{j,h}^+ < t_{\text{low}}$ and

$$t_{\text{up}} + \frac{1}{\tau_{\tilde{m}_{\text{up}}+1}^+(t_{\text{up}})} \geq t_{\text{low}}.$$

- d) Compute

$$\omega_{\tilde{m}_{\text{low}}-k+1,h}^- = t_{\text{low}} + \frac{1}{\tau_k^-(t_{\text{low}})} \quad \text{for } k = 1, \dots, \tilde{m}_{\text{low}}$$

where \tilde{m}_{low} is such that all $\omega_{\tilde{m}_{\text{low}}-k+1,h}^- > t_{\text{up}}$ and

$$t_{\text{low}} + \frac{1}{\tau_{\tilde{m}_{\text{low}}+1}^-(t_{\text{low}})} \leq t_{\text{up}}.$$

- e) If $\tilde{m}_{\text{low}} \neq \tilde{m}_{\text{up}}$, decrease h , set new $\mathcal{L} = \tilde{\mathcal{L}}_h \in \mathcal{F}$ and go back to c).
 Otherwise set $\tilde{m} = \tilde{m}_{\text{low}} = \tilde{m}_{\text{up}}$, decrease h , set new $\mathcal{L} = \tilde{\mathcal{L}}_h \in \mathcal{F}$ and continue from b).
 f) Exit with $\tilde{m}(\delta) = \tilde{m}$ and $\omega_j^\pm(\delta) = \omega_{j,h}^\pm$ for $j = 1, \dots, \tilde{m}$.

Let

$$(t_{\text{up}}, t_{\text{low}}) \cap \sigma(\mathcal{M}) = \{\omega_{k+1}, \dots, \omega_{k+m}\}$$

where

$$m = \text{Tr } \mathbf{1}_{(t_{\text{up}}, t_{\text{low}})}(\mathcal{M}) \quad \text{and} \quad k \geq 0.$$

Observe that, a priori, an interval (ω_j^-, ω_j^+) obtained as the output of Procedure 1 is not guaranteed to have a non-empty intersection with the spectrum of \mathcal{M} or in fact include precisely the eigenvalue ω_{k+j} . However, as it is established by the following lemma, the latter is certainly true for δ small enough.

Lemma 16. *There exist $t^0 > 0$ and $\delta_0 > 0$, ensuring the following for all $t_{\text{low}} \geq t^0$ and $\delta < \delta_0$.*

- a) *The conditional loop in Procedure 1 always exits in the regime $h \rightarrow 0$.*
 b) *$m(\delta) = m$.*
 c) *$\omega_j^-(\delta) \leq \omega_{k+j} \leq \omega_j^+(\delta)$ for all $j = 1, \dots, n$.*

Proof. Since $\nu_j^+(t_{\text{up}}) = \omega_{k+j} = \nu_{n-j+1}^-(t_{\text{low}})$ for all $j = 1, \dots, n$, Theorem 14 alongside with the assumption on \mathcal{F} , confirms the existence of $\omega_{j,h}^\pm$ in Procedure 1-c) and d), for all $j = 1, \dots, n$ whenever h is small enough. Moreover

$$\omega_{j,h}^+ \downarrow \omega_{k+j} \quad \text{and} \quad \omega_{j,h}^- \uparrow \omega_{k+j} \quad \text{as } h \rightarrow 0.$$

This ensures the validity of the lemma. \square

If the eigenfunctions of \mathcal{M} lie in $\mathcal{H}^{r+1}(\Omega)^6$, where r is the degree of the polynomials in (42), then

$$\omega_{j,h}^+ - \omega_{j,h}^- = O(h^{2r}).$$

This means that the exit rate of the conditional loop in Procedure 1 is also $O(h^{2r})$ as $h \rightarrow 0$.

A close examination of the constants involved in the proof of Theorem 14, indicates that they are of order $|t - \nu_1^\pm(t)|^{-1}$. See Theorem 10 and Corollary 11. Table 1 and other various numerical experiments not included in Section 7, strongly suggest that the accuracy improves significantly, as $t_{\text{up}} \downarrow \nu_1^-(t_{\text{up}})$ and $t_{\text{low}} \uparrow \nu_1^+(t_{\text{low}})$.

7. COMPUTATIONAL EXAMPLES

We now illustrate the practical applicability of the ideas discussed above by means of several examples. Two canonical references for benchmarks on the Maxwell eigenvalue problem are [15] and [9]. We validate some of the numerical bounds found below against these benchmarks. All the experiments presented are performed for $\epsilon = \mu = 1$ and some of them consider the so-called two-dimensional Maxwell problem.

If the domain Ω has a cylindrical symmetry, say $\Omega = \tilde{\Omega} \times (0, \pi)$ for $\tilde{\Omega} \subset \mathbb{R}^2$ open and sufficiently regular, then (1) decouples. Indeed, by performing a separation of variables, a non-zero ω is an eigenvalue of \mathcal{M}_1 if and only if either $\omega^2 = \lambda^2$ or $\omega^2 = \nu^2 + \rho^2$, where λ^2 is a Dirichlet eigenvalue of the Laplacian in $\tilde{\Omega}$, ν^2 is

a non-zero Neumann eigenvalue of the Laplacian in $\tilde{\Omega}$ and $\rho \in \mathbb{N}$. In turns the Neumann problem can be re-written as

$$(50) \quad \begin{cases} \operatorname{curl} \mathbf{E} = i\mu H & \text{in } \tilde{\Omega} \\ \operatorname{curl} H = -i\omega \mathbf{E} & \text{in } \tilde{\Omega} \\ \mathbf{E} \cdot \mathbf{t} = \mathbf{0} & \text{on } \partial\tilde{\Omega}, \end{cases}$$

for non-zero $(\mathbf{E}, H)^t \in L^2(\tilde{\Omega})^3$ and $\nu = \omega \in \mathbb{R}$, where

$$\mathbf{E} = \begin{pmatrix} E_1 \\ E_2 \end{pmatrix}, \quad \operatorname{curl} \mathbf{E} = \partial_x E_2 - \partial_y E_1, \quad \operatorname{curl} H = \begin{pmatrix} \partial_y H \\ -\partial_x H \end{pmatrix}$$

and \mathbf{t} is the unit tangent to $\partial\tilde{\Omega}$. This two-dimensional Maxwell problem suffers from all the complications concerning spectral pollution, as its three-dimensional counterpart.

We denote by $\tilde{\mathcal{M}}$ the self-adjoint operator associated to (50). This operator can be employed for numerical tests which can then be validated against numerical calculations for the original Neumann Laplacian via the Galerkin method, [15]. Indeed, note that the latter is a semi-definite operator with a compact resolvent, so it does not exhibit spectral pollution.

The ideas developed in Section 5 for the operator \mathcal{M} have analogues for $\tilde{\mathcal{M}}$. In the lower-dimensional examples presented below, we have chosen the finite element spaces on a corresponding triangulation \mathcal{T}_h of $\tilde{\Omega}$ as

$$\begin{aligned} \mathbf{V}_h^{r,k} &= \{\mathbf{v}_h \in [C^0(\tilde{\Omega})]^k : \mathbf{v}_h|_K \in [\mathbb{P}_r(K)]^k \forall K \in \mathcal{T}_h\} \quad (k = 1, 2) \\ \mathbf{V}_{h,0}^{r,2} &= \{\mathbf{v}_h \in \mathbf{V}_h^{r,2} : \mathbf{v}_h \times \mathbf{n} = \mathbf{0} \text{ on } \partial\tilde{\Omega}\} \quad \text{and} \\ \mathcal{L}_h &= \mathbf{V}_{h,0}^{r,2} \times \mathbf{V}_h^{r,1}. \end{aligned}$$

This ensures that $\mathcal{L}_h \subset \mathbf{D}(\tilde{\mathcal{M}})$.

7.1. Convex domains. For a convex domain, the eigenfunctions of (1) or (50) possess interior regularity. This leads to an improvement in accuracy as a consequence of (48). In this, the best possible case scenario, the Zimmermann-Mertins method for the resonant cavity problem achieves an optimal order of convergence in the context of the finite element method.

Accuracy of the enclosures on a square. In this set of experiments we consider $\tilde{\Omega} = \Omega_{\text{sqr}} = (0, \pi)^2 \subset \mathbb{R}^2$. The eigenvalues of $\tilde{\mathcal{M}}$ are $\omega = \pm\sqrt{l^2 + m^2}$ for $l, m \in \mathbb{N} \cup \{0\}$. In order to estimate ω_k^\pm numerically, we have picked

$$t_{\text{up}} = \frac{1}{4}\omega_{k-1} + \frac{3}{4}\omega_k \quad \text{and} \quad t_{\text{low}} = \frac{3}{4}\omega_k + \frac{1}{4}\omega_{k+1}$$

to machine precision. Here and below we substitute from the notation in previous sections the index j for eigenvalues by an index k , in order to highlight the fact that we do not always count multiplicities.

In our first experiment we have computed the enclosure width $\omega_k^+ - \omega_k^-$ for $k = 1, \dots, 100$ and $r = 1, 3, 5$. We have chosen $h = h(r)$ such that the subspaces \mathcal{L}_h have roughly the same dimension $\approx 61\text{K}$. Figure 1 shows the outcomes of this experiment. We have excluded enclosures with size above 10^{-1} . As it is natural to expect, for a fixed subspace \mathcal{L}_h , the accuracy deteriorates as the eigenvalue counting number j increases: high energy eigenfunctions present more oscillations

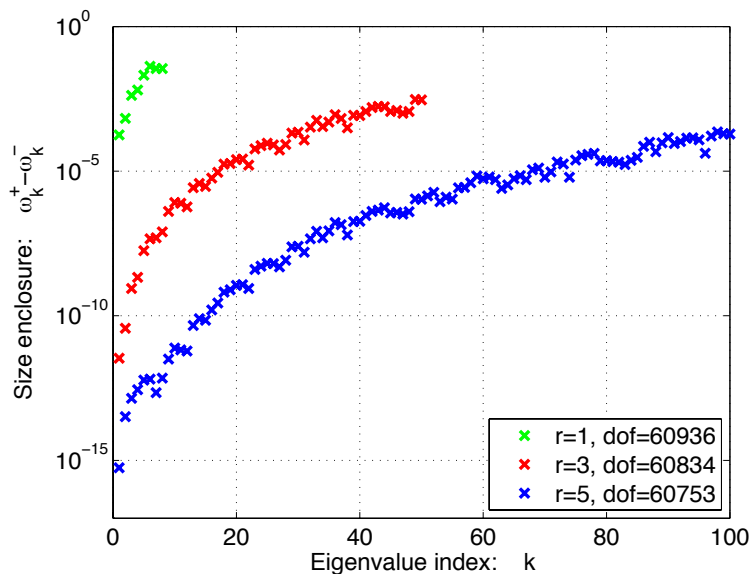


FIGURE 1. Semi-log graph associated to Ω_{sqr} . Vertical axis: $\omega_k^+ - \omega_k^-$. Horizontal axis: eigenvalue index k (not counting multiplicity). Here we use elements of order $r = 1, 3, 5$ on unstructured uniform meshes rendering roughly the same degrees of freedom.

that make their approximation more challenging. The accuracy increases with the polynomial order. The first 100 eigenvalues are approximated fairly accurately (note that $\omega_{100} = \sqrt{261}$) with polynomial order $r = 5$.

Convergence for a cube. We now consider numerical approximation of the eigenvalues of the three dimensional problem (1) for $\Omega = \Omega_{\text{cbe}} = (0, \pi)^3 \subset \mathbb{R}^3$. The non-zero eigenvalues are now $\omega = \pm\sqrt{l^2 + m^2 + n^2}$. The corresponding eigenfunctions are

$$\mathbf{E}(x, y, z) = \begin{pmatrix} \alpha_1 \cos(lx) \sin(my) \sin(nz) \\ \alpha_2 \sin(lx) \cos(my) \sin(nz) \\ \alpha_3 \sin(lx) \sin(my) \cos(nz) \end{pmatrix} \quad \forall \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix} \cdot \begin{pmatrix} l \\ m \\ n \end{pmatrix} = 0.$$

Here $\{l, m, n\} \subset \mathbb{N} \cup \{0\}$ and not two indices are allowed to vanish simultaneously. The vector $\underline{\alpha}$ determines the multiplicity of the eigenvalue for a given triplet (l, m, n) . That is, for example, $\omega = \sqrt{2}$ (the first positive eigenvalue) has multiplicity 3 corresponding to indices $\{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}$ each one of them contributing to one of the dimensions of the eigenspace. However, $\omega = \sqrt{3}$ (the second positive eigenvalue) corresponding to index $\{(1, 1, 1)\}$ has multiplicity 2 determined by $\underline{\alpha}$ on a plane.

In Figure 3 we have depicted the decrease in the enclosure width for the computation of the eigenvalue $\omega_2 = \sqrt{2}$ for Lagrange elements of order $r = 1, 2, 3$. We have chosen a sequence of unstructured tetrahedral mesh. The computed values for the slopes of the straight lines indicate that the enclosures obey the estimate

$$(51) \quad |\omega_j^\pm - \omega_j| \leq ch^{2r}.$$

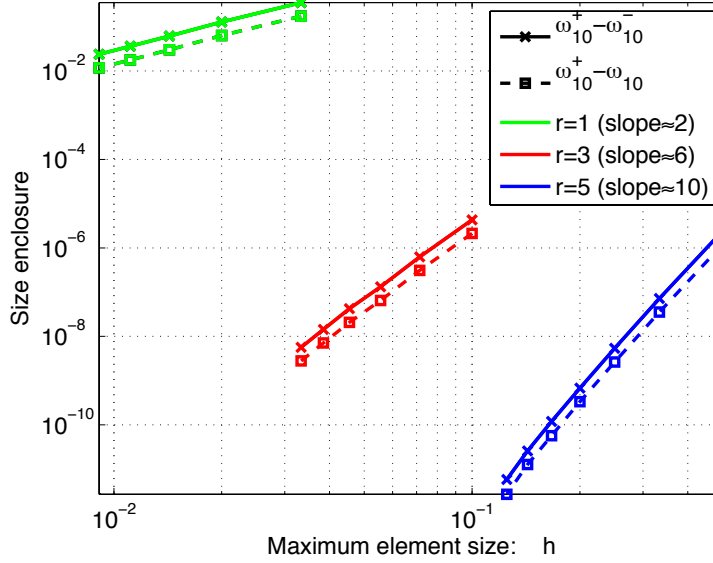


FIGURE 2. Log-log graph associated to Ω_{sqr} and $\omega_{10} = \sqrt{17}$. Vertical axis: enclosure width. Horizontal axis: Maximum element size h . Here we have chosen Lagrange elements of order $r = 1, 3, 5$ on a sequence of unstructured meshes.

Therefore the conclusion (48) of Theorem 14 will be sharp. Note that in the picture, we have considered both the exact residual and the length of the enclosure.

The slashed cube. We now assume that $\Omega = \Omega_{\text{sla}} = (0, \pi)^3 \setminus T \subset \mathbb{R}^3$. Here T is the closed tetrahedron with vertices $(0, 0, 0)$, $(\pi/2, 0, 0)$, $(0, \pi/2, 0)$ and $(0, 0, \pi/2)$. This domain does not have symmetries allowing a reduction into two-dimensions. However, as Ω_{sla} is fairly close to Ω_{cbe} , we should expect that the structure of the spectrum in the two cases is reminiscent of one another.

In our first experiment on this region, we compute benchmark eigenvalue enclosures for (1). The table to the right of Figure 4 shows the outcomes of implementing the Procedure 1. We have run an algorithm based on this procedure for each of three fixed choices of t_{up} and t_{low} (third and fourth columns) with $\delta = 10^{-2}$. We have picked the family of mesh so that no more than five iterations were required to achieve the needed accuracy. The parameter l in this table counts the number of eigenvalues to the right of t_{up} or to the left of t_{low} , respectively.

In this experiment we have chosen trial spaces made out of Lagrange elements of order $r = 3$. All the final eigenvalue enclosures have a length of at most 2×10^{-3} . The mesh used in the last iteration is depicted on the left of Figure 4.

From the table it seems clear that there is a cluster of eigenvalues at the bottom of the positive spectrum near $\sqrt{2}$. The latter is the first positive eigenvalue of Ω_{cbe} which is of multiplicity 3. It appears that this eigenvalue splits into a single eigenvalue at the bottom of the spectrum and a seemingly double eigenvalue slightly above it. Another cluster occurs at ω_4 and ω_5 with strong indication that this is a double eigenvalue. This pair is near $\sqrt{3}$, the second eigenvalue of Ω_{cbe} which is

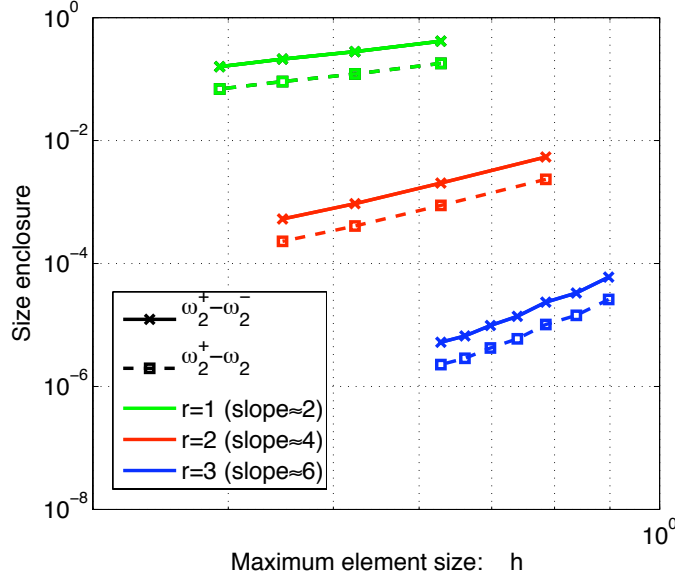
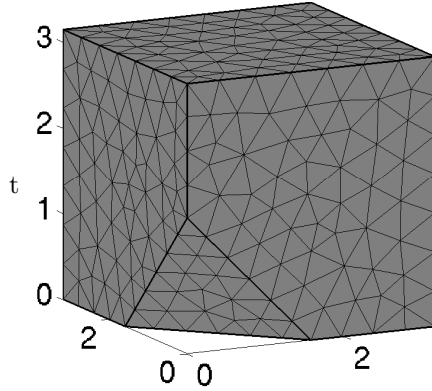


FIGURE 3. Log-log graph associated to Ω_{cbe} and $\omega_2 = \sqrt{2}$. Vertical axis: enclosure width. Horizontal axis: maximum element size h . Here we have chosen Lagrange elements of order $r = 1, 2, 3$ on a sequence of unstructured meshes.



k	$(\omega_k)_{-}^{+}$	$t_{up}(l)$	$t_{low}(l)$
1	1.412_{000}^{236}	0.5 (1)	1.6 (3)
2	1.430_{560}^{672}	0.5 (2)	1.6 (2)
3	1.430_{577}^{673}	0.5 (3)	1.6 (1)
4	1.755_{043}^{308}	1.5 (1)	2.1 (2)
5	1.755_{063}^{329}	1.5 (2)	2.1 (1)
6	2.22_{053}^{200}	1.8 (1)	2.6 (5)
7	2.237_{434}^{667}	1.8 (2)	2.6 (4)
8	2.237_{459}^{684}	1.8 (3)	2.6 (3)
9	2.239_{387}^{533}	1.8 (4)	2.6 (2)
10	2.270_{558}^{778}	1.8 (5)	2.6 (1)

FIGURE 4. Benchmark spectral approximation for Ω_{sla} . In the table we compute interval of enclosure for the first 10 eigenvalues of (1). In order to obtain this calculation we have employed Procedure 1. The trial spaces are made of Lagrange elements of order $r = 3$. The final mesh is the one shown on the right side. Total number of DOF=117102.

indeed double. The next eigenvalues for the cube are 2 and $\sqrt{5}$ with total multiplicity 5. It is natural to conjecture that ω_j for $j = 5, \dots, 10$ are perturbations of these eigenvalues, but the data shown in the table is inconclusive.

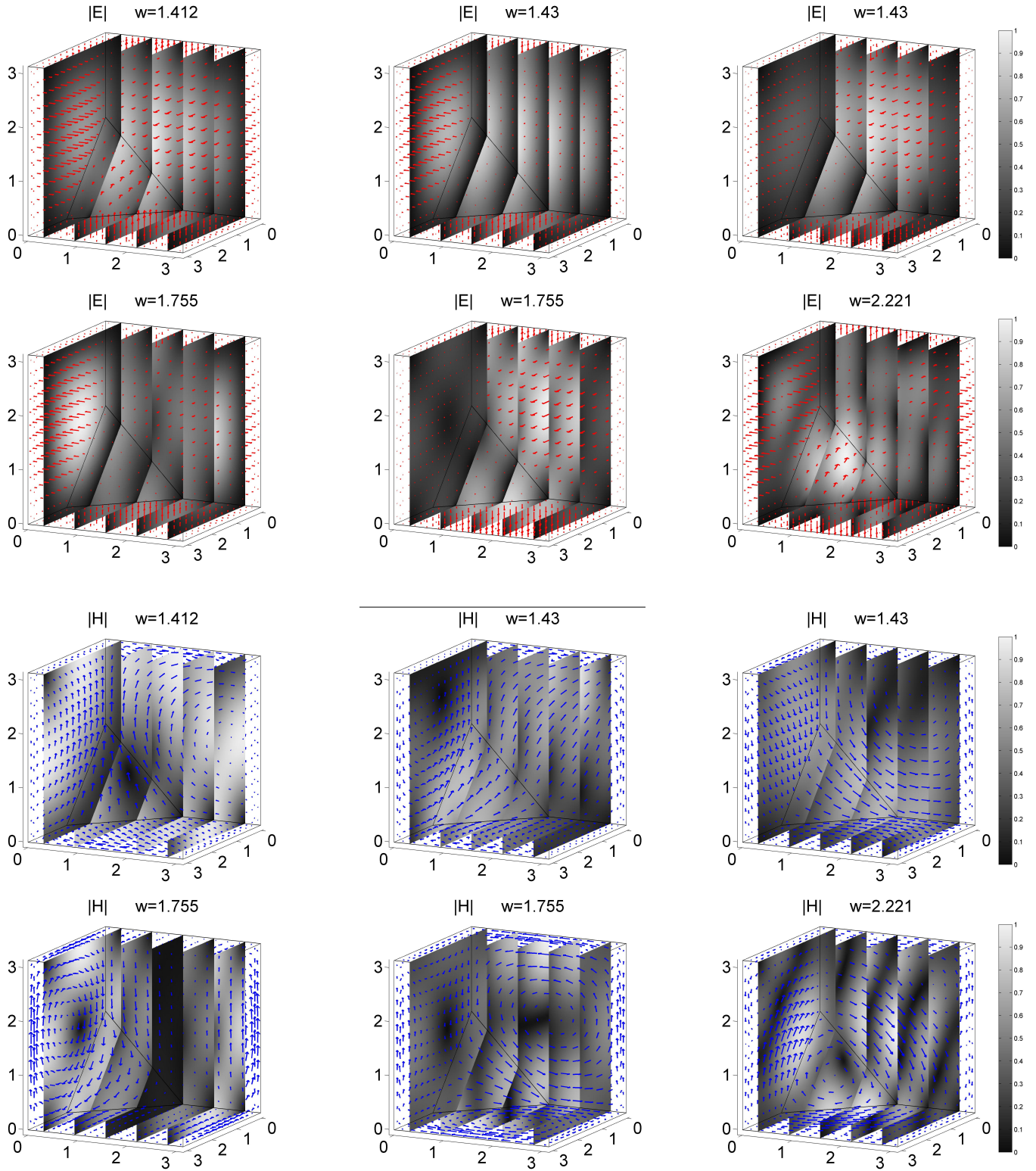


FIGURE 5. The first six eigenfunctions on Ω_{sla} for the first six positive eigenvalues. Densities $|E|$ (top) and $|H|$ (bottom). Corresponding arrow fields E (red) and H (blue) on $\partial\Omega_{\text{sla}}$.

For our next experiment on this region, we have estimated numerically the electromagnetic fields corresponding to index up to 6 from the table in Figure 4. The purpose of the experiment is to set benchmarks for the eigenfunctions on Ω_{sla} and simultaneously illustrate Theorem 15. In Figure 5 we depict the density of electric and magnetic fields, $|\mathbf{E}|$ and $|\mathbf{H}|$ both normalized to having maximum equal to 1. We also show arrows pointing towards the direction of these fields on $\partial\Omega_{\text{sla}}$.

The mesh employed for these calculations is the one shown in Figure 4. According to Theorem 15 and the data presented in the table, the shown eigenfunctions should be close to the exact eigenfunctions in the curl norm. We remark that for both experiments on Ω_{sla} a reasonable accuracy has been achieved even for the fairly coarse mesh depicted in Figure 4.

7.2. Non-convex domains. The numerical approximation of the eigenfrequencies and electromagnetic fields in the resonant cavity is known to be challenging when the domain is not convex. The main reason for this is the fact that the electromagnetic field might have a singularity and a low degree of regularity at re-entrant corners. See for example the discussion after [22, Lemma 3.56] and references therein.

In some of the examples of this section we consider a mesh adapted to the geometry of Ω . However, we do not pursue any specialized mesh refinement strategy. We show below that, even in the case where there is poor approximation due to low regularity of the eigenspace, the method presented in this paper provides valuable information about the localization of the eigenvalues of (1).

The L-shaped domain. The region $\tilde{\Omega} = \Omega_L = (0, \pi)^2 \setminus [0, \pi/2]^2$ is a classical benchmark domain both for the Maxwell and the Helmholtz problems, and it has been extensively examined in the past. Numerical computations for the eigenvalues of \mathcal{M} were reported in [9, Table 5] via an implementation based on a mixed formulation of (50) and the help of edge finite elements. See also [15]. We now show how to achieve accurate enclosures for these eigenvalues with the help of nodal finite elements.

For this next set of experiments we consider unstructured triangulations of the domain, refined around the re-entrant corner. The polynomial order is set to $r = 3$. Figures 6 - 8 compile our findings.

We produced the table in Figures 6 by implementing Procedure 1 in the same fashion as for the case of Ω_{sla} discussed previously. For comparison in the second column of this table we have included the benchmark eigenvalue estimations found in [9] and [15]. Note that some of the computed eigenvalues associated to the mixed formulation are lower bounds of the true eigenvalues, and some, like the 9th, are upper bounds. This confirms that the latter approach is in general un-hierarchical as previously suggested in the literature.

From the third column of the table, it is clear that the accuracy depends on the regularity of the corresponding eigenspaces. The eigenfunctions associated to $\omega_3 = \omega_4 = 2$ and $\omega_7 = \sqrt{8}$ are found by gluing together corresponding eigenfunctions of (1) on squares of side $\pi/2$. These eigenfunctions are smooth in the interior of Ω_L , while those associated to ω_1 and ω_2 present singularities around the re-entrant corner. The electric field component of the former is known to be outside $\mathcal{H}^1(\Omega_L)^2$ while that of the latter is in $\mathcal{H}^1(\Omega_L)^2$. This explains the significant gain in accuracy in the calculation of ω_2 with respect to the one of ω_1 .

k	ω_j from [9] (from [15])	$(\omega_j)_\pm^+$	$t_{\text{up}}(l)$	$t_{\text{low}}(l)$
1	0.768192684 (0.773334985176)	$0.773334\frac{991}{694}$	0.1 (1)	2.1 (4)
2	1.196779010 (1.19678275574)	$1.1967827557\frac{026}{761}$	0.1 (2)	2.1 (3)
3	1.999784988 (2.00000000000)	$\frac{2.00000000064}{1.99999999933}$	1.5 (1)	2.5 (4)
4	1.999784988 (2.00000000000)	$\frac{2.00000000067}{1.99999999936}$	1.5 (2)	2.5 (3)
5	2.148306309 (2.14848368266)	$2.14848368\frac{365}{199}$	3.1 (5)	1.5 (3)
6	2.252760528	$2.25729\frac{896}{776}$	1.5 (4)	3.1 (4)
7	2.828075317	$2.8284271\frac{354}{186}$	1.5 (5)	3.7 (4)
8	2.938491109	$2.94671\frac{343}{112}$	1.5 (6)	3.7 (3)
9	3.075901493	$3.0758929\frac{738}{571}$	1.5 (7)	3.7 (2)
10	3.390427701	$3.3980\frac{724}{676}$	1.5 (8)	3.7 (1)

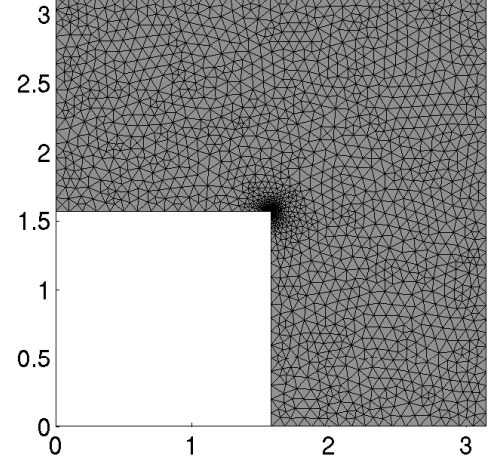


FIGURE 6. Enclosures for the first 10 positive eigenvalues of \mathcal{M} on the region Ω_L . The next eigenvalue is above 3.7. Here Procedure 1 has been implemented on Lagrange elements of order 3. The final mesh shown on the right has a number of DOF=56055. The mesh has a maximum element size $h = 0.1$ and has been refined at $(\pi/2, \pi/2)$. For comparison on the second column we include the eigenvalue estimations found in [9] and [15].

Figure 7 depicts in log-log scale residuals versus maximum element size. We have considered here Lagrange elements of order $r = 3$ and $r = 5$. The hierarchy of mesh (not shown) was chosen unstructured, but with an uniform distribution of nodes. Since the eigenfunctions associated to ω_1 and ω_2 have a limited regularity, then there is no noticeable improvement of convergence order as r increases. As the third eigenfunction is smooth, it does obey the estimate (51).

Benchmark approximated eigenfunctions are depicted in Figure 8. The mesh employed to produce these graphs is the one shown on the right of Figure 6. As some of the electric fields have a singularity at $(\pi/2, \pi/2)$ we have normalized each individual plot to a range in the interval $[0, 1]$.

The Fichera domain. In this next experiment we approximate for the eigenpairs of (1) associated to the region $\Omega = \Omega_F = (0, \pi)^3 \setminus [0, \pi/2]^3 \subset \mathbb{R}^3$ numerically.

Some of the eigenvalues can be obtained by domain decomposition and the corresponding eigenfunctions are regular. For example, eigenfunctions on the cube of side $\pi/2$ can be assembled in the obvious fashion, in order to build eigenfunctions on Ω_F . Therefore the set $\{\pm 2\sqrt{l^2 + m^2 + n^2}\}$ where not two indices vanish simultaneously certainly lie in $\sigma(\mathcal{M})$. The first eigenvalue in this set is $2\sqrt{2}$. We conjecture that there are exactly 15 eigenvalues in the interval $(0, 2\sqrt{2})$. Furthermore, we conjecture that the multiplicity counting of the spectrum in this interval is

$$1, 2, 3, 2, 1, 2, 1, 3.$$

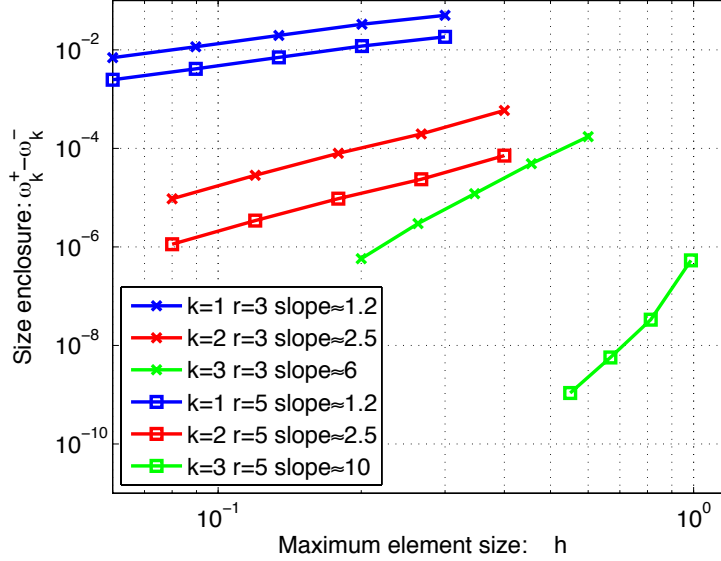


FIGURE 7. Compared order of approximation for different eigenvalues in the region Ω_L . The log-log plot shows residual versus maximum element size h for the calculation of enclosures for ω_k where $k = 1, 2, 3$ and \mathcal{L} is generated by Lagrange elements of order $r = 3$ and $r = 5$. Note that $(\mathbf{E}, H) \notin \mathcal{H}^s(\Omega_L)^3$ for $k = 1$ and $s = 1$, and for $k = 2$ and $s = 1.5$. On the other hand, for $k = 3$ we have (\mathbf{E}, H) smooth, as the eigenfunction is also solution of (1) on a square of side $\pi/2$.

The table on the right of Figure 9 shows a numerical estimation of these eigenvalues. Here we have fixed $t_{\text{up}} = 0.2$ and $t_{\text{low}} = 2.8$. We have considered a family of mesh refined along the re-entrant edges. The final mesh is shown on the left side of Figure 9. We have stopped the algorithm when the tolerance $\delta = 0.05$ has been achieved. However, note that the accuracy is much higher for the indices $k = 2, 3, 9, 10, 11, 15$.

The slight numerical discrepancy shown in the table for the seemingly multiple eigenvalues appears to be a consequence of the fact that the meshes employed are not entirely symmetric with respect to permutation of the spacial coordinates. Figure 10 includes the corresponding approximated eigenfunctions. The mesh employed for this calculation is the same as that of Figure 9.

The slit square. As mentioned in Section 6, for a single trial space \mathcal{L} , the accuracy of the eigenvalue bounds produced by the Zimmermann-Mertins method depends on the position of t relative to the spectrum of \mathcal{M} . In this final experiment we demonstrate that this dependence might vary significantly with t . The numerical evidence suggests that a good choice of t_{up} and t_{low} plays a role in the design of efficient algorithms for eigenvalue calculation based on this method.

Let $\tilde{\Omega} = (0, \pi)^2 \setminus S$ for $S = [\pi/2, \pi] \times \{\pi/2\}$. Benchmarks on the eigenvalues of (50) are known by means of solving numerically the corresponding Neumann

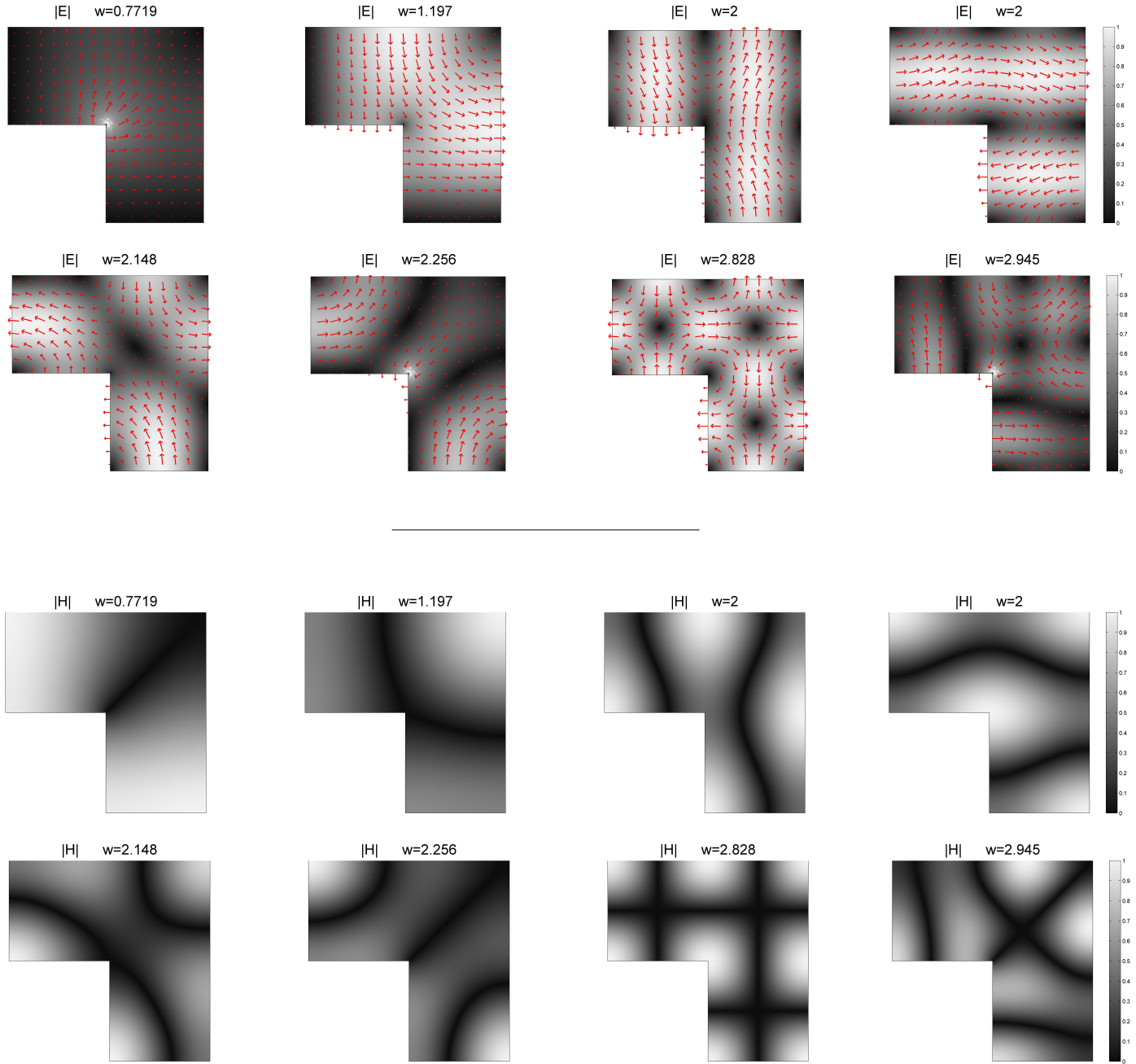


FIGURE 8. Eigenfunctions on Ω_L associated to the first eight positive eigenvalues. Densities $|E|$ (top) and $|H|$ (bottom). Corresponding arrow fields \mathbf{E} . We have normalized each individual density to have as maximum the value 1.

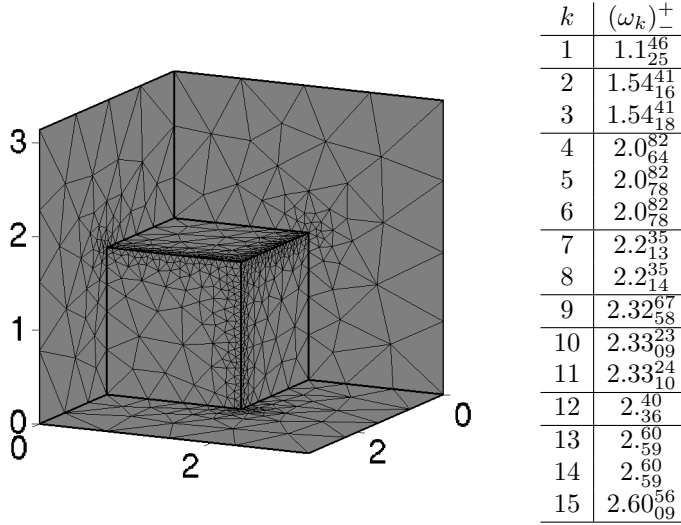


FIGURE 9. Spectral enclosures for the spectrum lying on the interval $(0, 2\sqrt{2})$ for the Fichera domain Ω_F . Here we have fixed $t_{\text{up}} = 0.2$ and $t_{\text{low}} = 2.8$. We considered mesh refined at the re-entrant edges as shown on the left. The final number of DOF=208680.

RF	DOF	$t_{\text{low}} = 1.95$ $(j = 1) \omega_3^-$	$t_{\text{low}} = 2.05$ $(j = 3) \omega_3^-$	$t_{\text{up}} = 1.05$ $(j = 1) \omega_3^+$	$t_{\text{up}} = 0.7$ $(j = 3) \omega_3^+$
1	4143	1.24764	1.26640	1.50395	1.3436
0.1	9648	1.25029	1.26830	1.49282	1.3336
0.01	74226	1.25063	1.26846	1.48899	1.3274

TABLE 1. Numerical experiment showing the dependence of the accuracy in the Zimmermann-Mertins method on the choice of t_{up} and t_{low} . It is preferable to pick t_{up} and t_{low} as far as possible from ω than to increase the dimension of \mathcal{L} .

Laplacian problems, [15]. The first seven positive eigenvalues are

$$\omega_1 \approx 0.647375015, \omega_2 = 1, \omega_3 \approx 1.280686161, \\ \omega_4 = \omega_5 = 2, \omega_6 \approx 2.096486081 \quad \text{and} \quad \omega_7 \approx 2.229523505.$$

The eigenfunctions associated to ω_2, ω_4 and ω_5 are smooth, as they are also eigenfunctions on Ω_{sqr} . On the other hand, ω_1 and ω_3 , correspond to singular eigenfunctions. Standard nodal elements are completely unsuitable for the computation of these eigenvalues, even with a significant refinement of the mesh on S .

Table 1 shows computation of ω_3^{\pm} on a mesh that is increasingly refined at S with a factor RF for two pairs of choices of t_{up} and t_{low} . Here $h = 0.1$ and we consider Lagrange elements of order $r = 1$. The choice of t_{up} and t_{low} further from ω_3 even with the very coarse mesh, provide qualitatively sharper ω_3^{\pm} than the other choices even with the finer mesh.

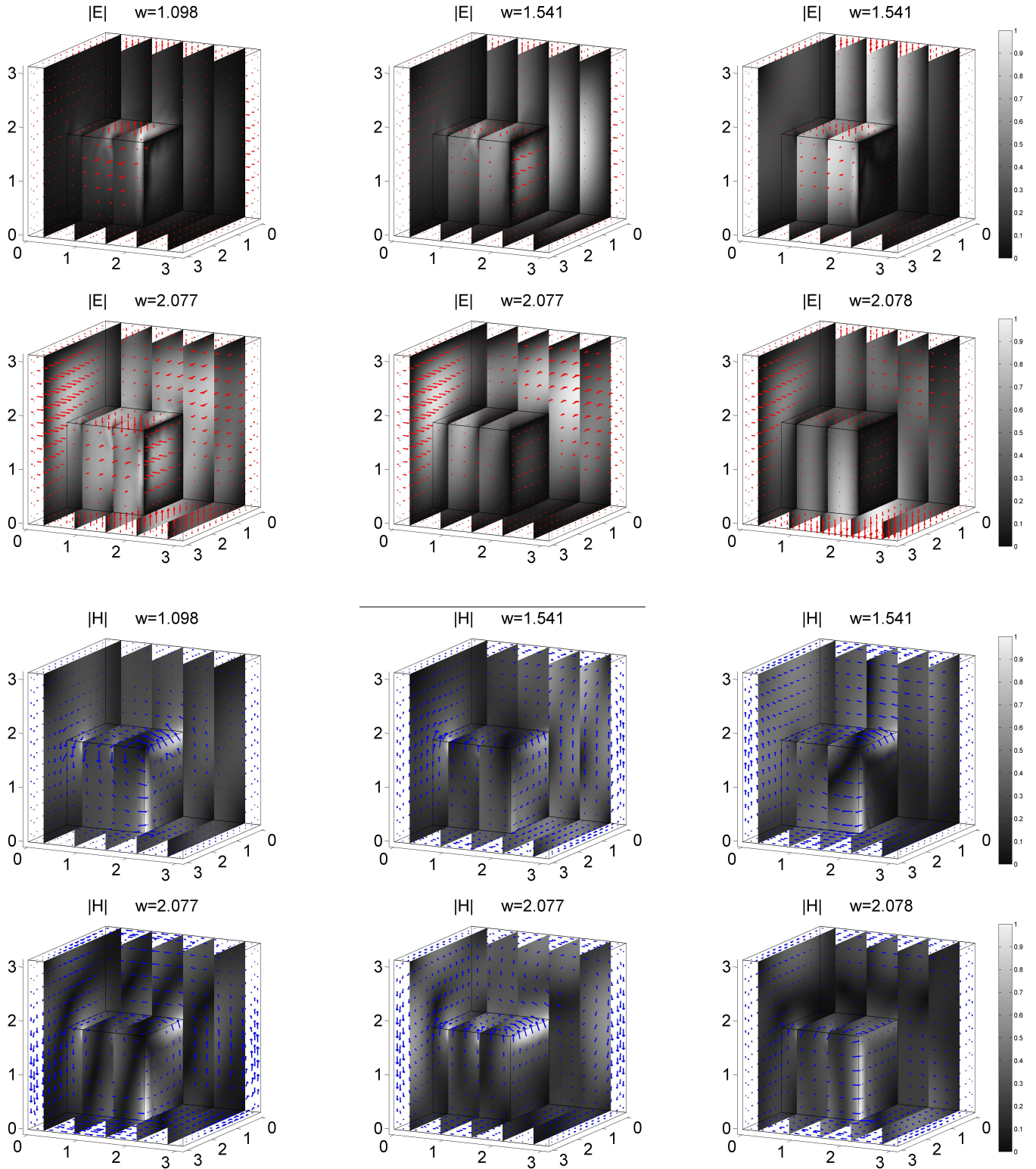


FIGURE 10. The first six eigenfunctions on Ω_F for the first six positive eigenvalues. Densities $|E|$ (top) and $|H|$ (bottom). Corresponding arrow fields E (red) and H (blue) on $\partial\Omega_F$.

APPENDIX A. FURTHER GEOMETRICAL PROPERTIES OF $F_{\mathcal{L}}^j(t)$

Various extensions of Lemma 2 to the case $j > 1$ are possible, however it is difficult to write these results in a neat fashion. The proposition below is one such an extension.

The following generalization of Danskin's Theorem is a direct consequence of [6, Theorem D1]. Let $J \subset \mathbb{R}$ be an open segment. Denote by

$$\partial_t^\pm f(t) = \lim_{\tau \rightarrow 0^+} \pm \frac{f(t \pm \tau) - f(t)}{\tau},$$

the one-side derivatives of a function $f : J \rightarrow \mathbb{R}$. Let \mathcal{V} be a compact topological. For given $\mathcal{J} : J \times \mathcal{V} \rightarrow \mathbb{R}$ we write

$$\tilde{\mathcal{J}}(t) = \max_{v \in \mathcal{V}} \mathcal{J}(t, v) \quad \text{and} \quad \tilde{\mathcal{V}}(t) = \left\{ \tilde{v} \in \mathcal{V} : \tilde{\mathcal{J}}(t) = \mathcal{J}(t, \tilde{v}) \right\}.$$

Lemma 17. *If the map \mathcal{J} is upper semi-continuous and $\partial_t^\pm \mathcal{J}(t, v)$ exist for all $(t, v) \in J \times \mathcal{V}$, then also $\partial_t^\pm \tilde{\mathcal{J}}(t)$ exist for all $t \in J$ and*

$$(52) \quad \partial_t^\pm \tilde{\mathcal{J}}(t) = \max_{\tilde{v} \in \tilde{\mathcal{V}}(t)} \partial_t^\pm \mathcal{J}(t, \tilde{v}).$$

In the statement of this lemma, note that the left and right derivatives of both \mathcal{J} and $\tilde{\mathcal{J}}$ might possibly be different.

Proposition 18. *Let $j = 1, \dots, n$ and $t \in \mathbb{R}$ be fixed. The following assertions are equivalent.*

- a) $|F_{\mathcal{L}}^j(t) - F_{\mathcal{L}}^j(s)| = |t - s|$ for some $s \neq t$.
- b) There exists an open segment $J \subset \mathbb{R}$ containing t in its closure, such that

$$|F_{\mathcal{L}}^j(t) - F_{\mathcal{L}}^j(s)| = |t - s| \quad \forall s \in \bar{J}.$$

- c) There exists an open segment $J \subset \mathbb{R}$ containing t in its closure, such that

$$\forall s \in J, \text{ either } \mathcal{L} \cap \mathcal{E}_{s+F_{\mathcal{L}}^j(s)} \neq \{0\} \quad \text{or} \quad \mathcal{L} \cap \mathcal{E}_{s-F_{\mathcal{L}}^j(s)}(A) \neq \{0\}.$$

Proof.

a) \Rightarrow b). Assume a). Since $r \mapsto r \pm F_{\mathcal{L}}^j(r)$ are continuous and monotonically increasing, then they have to be constant in the closure of

$$J = \{\tau t + (1 - \tau)s : 0 < \tau < 1\}.$$

This is precisely b).

b) \Rightarrow c). Assume b). Then $s \mapsto F_{\mathcal{L}}^j(s)$ is differentiable in J and its one-side derivatives are equal to 1 or -1 in the whole of this interval. For this part of the proof, we aim at applying (52), in order to get another expression for these derivatives.

Let \mathcal{F}_j be the family of $(j-1)$ -dimensional linear subspaces of \mathcal{L} . Identify an orthonormal basis of \mathcal{L} with the canonical basis of \mathbb{C}^n . Then any other orthonormal basis of \mathcal{L} is represented by a matrix in $O(n)$, the orthonormal group. By picking the first $(j-1)$ columns of these matrices, we cover all possible subspaces $V \in \mathcal{F}_j$. Indeed we just have to identify $(v_1 | \dots | v_{j-1})$ for $[v_{kl}]_{kl=1}^n \in O(n)$ with $V = \text{Span}\{v_k\}_{k=1}^{j-1}$.

Let

$$\mathcal{K}_j = \left\{ (v_1, \dots, v_{j-1}) : [v_{kl}]_{kl=1}^n \in O(n) \right\} \subset \underbrace{\mathbb{C}^n \times \dots \times \mathbb{C}^n}_{j-1}.$$

Then \mathcal{K}_j is a compact subset in the product topology of the right hand side. According to (11),

$$F_{\mathcal{L}}^j(s) = \max_{(\underline{v}_1, \dots, \underline{v}_{j-1}) \in \mathcal{K}_j} g(s; \underline{v}_1, \dots, \underline{v}_{j-1})$$

where

$$g(s; \underline{v}_1, \dots, \underline{v}_{j-1}) = \min_{\substack{(a_1, \dots, a_{j-1}) \in \mathbb{C}^{j-1} \\ \sum |a_k|^2 = 1}} \left| \sum a_k \tilde{v}_k \right|_s.$$

Here we have used the correspondence between $\underline{v}_k \in \mathbb{C}^n$ and $\tilde{v}_k \in \mathcal{L}$ in the orthonormal basis set above. We write

$$g(r, V) = g(r; \underline{v}_1, \dots, \underline{v}_{j-1}) \quad \text{for } V = \text{Span}\{\tilde{v}_k\}_{k=1}^{j-1} \in \mathcal{F}_j.$$

The map $g : J \times \mathcal{K}_j \rightarrow \mathbb{R}^+$ is the minimum of a differentiable function, so the hypotheses of Lemma 17 are satisfied by $\mathcal{J} = -g$. Hence, by virtue of (52),

$$\partial_s^\pm g(s, V) = \min_{\substack{u \in \mathcal{L} \ominus V, \|u\|=1 \\ |u|_s = g(s, V)}} \left(\frac{\text{Re } l_s(u, u)}{|u|_s} \right).$$

As minima of continuous functions, $g(s, V)$ and $\partial_s^\pm g(s, V)$ are upper semi-continuous. Therefore, a further application of Lemma 17 yields

$$\begin{aligned} \partial_s^\pm F_{\mathcal{L}}^j(s) &= \max_{\substack{(\underline{v}_1, \dots, \underline{v}_{j-1}) \in \mathcal{K}_j \\ g(s; \underline{v}_1, \dots, \underline{v}_{j-1}) = F_{\mathcal{L}}^j(s)}} \partial_s^\pm g(s, \underline{v}_1, \dots, \underline{v}_{j-1}) \\ &= \max_{\substack{V \in \mathcal{F}_j \\ g(s, V) = F_{\mathcal{L}}^j(s)}} \min_{\substack{u \in \mathcal{L} \ominus V, \|u\|=1 \\ |u|_s = g(s, V)}} \left(\frac{\text{Re } l_s(u, u)}{|u|_s} \right). \end{aligned}$$

Now, this shows that

$$\left| \max_{\substack{V \in \mathcal{F}_j \\ g(s, V) = F_{\mathcal{L}}^j(s)}} \min_{\substack{u \in \mathcal{L} \ominus V, \|u\|=1 \\ |u|_s = g(s, V)}} \left(\frac{\text{Re } l_s(u, u)}{|u|_s} \right) \right| = 1.$$

As \mathcal{L} is finite dimensional, there exists a vector $u \in \mathcal{L}$ satisfying $|u|_s = F_{\mathcal{L}}^j(s)$ such that

$$\frac{|\text{Re } l_s(u, u)|}{|u|_s} = 1.$$

Thus $|\text{Re} \langle (A - s)u, u \rangle| = \langle (A - s)u, (A - s)u \rangle = F_{\mathcal{L}}^j(s)$. Hence, according to the ‘‘equality’’ case in the Cauchy-Schwarz inequality, u must be an eigenvector of A associated with either $s + F_{\mathcal{L}}^j(s)$ or $s - F_{\mathcal{L}}^j(s)$. This is precisely c .

$c) \Rightarrow a)$. Under the condition c), there exists an open segment $\tilde{J} \subseteq J$, possibly smaller, such that $t \in \tilde{J}$ and $F_{\mathcal{L}}^j(s) = \mathfrak{d}_j(s)$ for all $s \in \tilde{J}$. As $|\mathfrak{d}_j(s) - \mathfrak{d}_j(r)| = |s - r|$, then either $a)$ is immediate, or it follows by taking $r \rightarrow t$. \square

APPENDIX B. A CORE COMSOL V4.3 LIVE LINK CODE

```
%
% Core Comsol V4.3 LiveLink code for computing
% fundamental frequencies on a resonant cavity
% with perfect conductivity conditions
% the test geometry below is the Fichera domain.
```

```

%
%
%       Gabriel Barrenchea, Lyonell Boulton
%       and Nabile Boussaid
%
%               November 2012
%

% INITIALIZATION OF THE MODEL FROM SCRATCHES

model = ModelUtil.create('Model');
geom1=model.geom.create('geom1', 3);
mesh1=model.mesh.create('mesh1', 'geom1');
w=model.physics.create('w', 'WeakFormPDE', 'geom1',
                      {'E1','E2', 'E3', 'H1', 'H2', 'H3'});

% CREATING THE GEOMETRY - IN THIS CASE THE FICHERA DOMAIN

hex1=geom1.feature.create('hex1', 'Hexahedron');
hex1.set('p',{ '0' '0' '0' '0' 'pi' 'pi' 'pi' 'pi';
              '0' '0' 'pi' 'pi' '0' '0' 'pi' 'pi';
              '0' 'pi' 'pi' '0' '0' 'pi' 'pi' '0'});
hex2=geom1.feature.create('hex2', 'Hexahedron');
hex2.set('p',{ '0' '0' '0' '0' 'pi/2' 'pi/2' 'pi/2' 'pi/2';
              '0' '0' 'pi/2' 'pi/2' '0' '0' 'pi/2' 'pi/2';
              '0' 'pi/2' 'pi/2' '0' '0' 'pi/2' 'pi/2' '0'});
dif1 = geom1.feature.create('dif1', 'Difference');
dif1.selection('input').set({'hex1'});
dif1.selection('input2').set({'hex2'});
geom1.run;

%CREATING THE GEOMETRY
model.mesh('mesh1').automatic(false);
model.mesh('mesh1').feature('size').set('custom', 'on');
model.mesh('mesh1').feature('size').set('hmax', '.8');
mesh1.run;

% PARAMETER t WHERE TO LOOK FOR EIGENVALUES
parat=2.2;

% WHETHER TO LOOK FOR THE EIGENVALUES TO THE LEFT (-) OR RIGHT (+) AND WHERE
ABOUT
shi=-.3;
model.param.set('tt', num2str(parat));
searchtau=shi;

% FINITE ELEMENTS TO USE AND ORDER

```



```

w.prop('ShapeProperty').set('shapeFunctionType', 'shlag');
w.prop('ShapeProperty').set('order', 3);

% PHYSICS
w.feature('wfeq1').set('weak',1 ,'(H3y-H2z)*(H3y_test-H2z_test)-
i*2*tt*(H3y-H2z)*E1_test+tt^2*E1*E1_test+(i*(H3y-H2z)-tt*E1)*E1t_test');
w.feature('wfeq1').set('weak',2 ,'(H1z-H3x)*(H1z_test-H3x_test)-
i*2*tt*(H1z-H3x)*E2_test+tt^2*E2*E2_test+(i*(H1z-H3x)-tt*E2)*E2t_test');
w.feature('wfeq1').set('weak',3 ,'(H2x-H1y)*(H2x_test-H1y_test)-
i*2*tt*(H2x-H1y)*E3_test+tt^2*E3*E3_test+(i*(H2x-H1y)-tt*E3)*E3t_test');
w.feature('wfeq1').set('weak',4 ,'(E3y-E2z)*(E3y_test-E2z_test)+
i*2*tt*(E3y-E2z)*H1_test+tt^2*H1*H1_test+((-i)*(E3y-E2z)-tt*H1)*H1t_test');
w.feature('wfeq1').set('weak',5 ,'(E1z-E3x)*(E1z_test-E3x_test)+
i*2*tt*(E1z-E3x)*H2_test+tt^2*H2*H2_test+((-i)*(E1z-E3x)-tt*H2)*H2t_test');
w.feature('wfeq1').set('weak',6 ,'(E2x-E1y)*(E2x_test-E1y_test)+
i*2*tt*(E2x-E1y)*H3_test+tt^2*H3*H3_test+((-i)*(E2x-E1y)-tt*H3)*H3t_test');

% BOUNDARY CONDITIONS
cons1=model.physics('w').feature.create('cons1', 'Constraint');
cons1.set('R', 2, 'E2');
cons1.set('R', 3, 'E3');
cons1.selection.set([1 8 9]);
cons2=model.physics('w').feature.create('cons2', 'Constraint');
cons2.set('R', 1, 'E1');
cons2.set('R', 3, 'E3');
cons2.selection.set([2 5 7]);
cons3=model.physics('w').feature.create('cons3', 'Constraint');
cons3.set('R', 1, 'E1');
cons3.set('R', 2, 'E2');
cons3.selection.set([3 4 6]);

% HOW MANY EIGENVALUES TO LOOK FOR AROUND t
neval=3;

% SOLVING THE MODEL
std1=model.study.create('std1');
model.study('std1').feature.create('eigv', 'Eigenvalue');
model.study('std1').feature('eigv').set('shift', num2str(searchtau));
model.study('std1').feature('eigv').set('neigs', neval);
std1.run;

% STORING SOLUTION FOR POST PROCESSING
[SZ,NDOFS,DATA,NAME,TYPE]= mphgetp(model,'solname','sol1');

% DISPLAYING SOLUTION
for inde=1:neval,
tauinv=(real(DATA(inde)));
bd=parat+tauinv;

```

```

if tauinv<0, disp(['lower= ',num2str(bd,10)]);
else disp(['upper= ',num2str(bd,10)]);
end
disp(['DOF= ',num2str(NDOFS)])
end

```

ACKNOWLEDGEMENTS

We kindly thank Michael Levitin and Stefan Neuwirth for their suggestions during the preparation of this manuscript. We kindly thank Université de Franche-Comté, University College London and the Isaac Newton Institute for Mathematical Sciences, for their hospitality. Funding was provided by MOPNET, the British-French project PHC Alliance (22817YA), the British Engineering and Physical Sciences Research Council (EP/I00761X/1) and the French Ministry of Research (ANR-10-BLAN-0101).

REFERENCES

- [1] C. AMROUCHE, C. BERNARDI, M. DAUGE, AND V. GIRAULT, Vector potentials in three-dimensional non-smooth domains, *Math. Methods Appl. Sci.*, 21 (1998), pp. 823–864.
- [2] D. ARNOLD, R. FALK, AND R. WINTHER, Finite element exterior calculus: from hodge theory to numerical stability, *Bulletin of the American Mathematical Society*, 47 (2010), pp. 281–354.
- [3] G. BARRENECHEA, L. BOULTON, AND N. BOUSSAÏD, Some remarks on the spectral properties of the maxwell operator on rough domains and domains with symmetries, in preparation.
- [4] H. BEHNKE, Lower and upper bounds for sloshing frequencies, *Inequalities and Applications*, (2009), pp. 13–22.
- [5] H. BEHNKE AND U. MERTINS, Bounds for eigenvalues with the use of finite elements, *Perspectives on Enclosure Methods*, (2001), p. 119.
- [6] P. BERNHARD AND A. RAPAPORT, On a theorem of Danskin with an application to a theorem of von Neumann-Sion, *Nonlinear Anal.*, 24 (1995), pp. 1163–1181.
- [7] M. BIRMAN AND M. SOLOMYAK, The self-adjoint Maxwell operator in arbitrary domains, *Leningrad Math. J.*, 1 (1990), pp. 99–115.
- [8] D. BOFFI, Finite element approximation of eigenvalue problems, *Acta Numer.*, 19 (2010), pp. 1–120.
- [9] D. BOFFI, P. FERNANDES, L. GASTALDI, AND I. PERUGIA, Computational models of electromagnetic resonators: analysis of edge element approximation, *SIAM J. Numer. Anal.*, 36 (1999), pp. 1264–1290 (electronic).
- [10] A. BONITO AND J.-L. GUERMOND, Approximation of the eigenvalue problem for the time harmonic Maxwell system by continuous Lagrange finite elements, *Math. Comp.*, 80 (2011), pp. 1887–1910.
- [11] L. BOULTON AND M. STRAUSS, Eigenvalue enclosures for the MHD operator, *BIT Numerical Mathematics*, (2012).
- [12] J. H. BRAMBLE, T. V. KOLEV, AND J. E. PASCIAK, The approximation of the Maxwell eigenvalue problem using a least-squares method, *Math. Comp.*, 74 (2005), pp. 1575–1598 (electronic).
- [13] A. BUFFA, P. CIARLET, AND E. JAMELOT, Solving electromagnetic eigenvalue problems in polyhedral domains, *Numer. Math.*, 113 (2009), pp. 497–518.
- [14] F. CHATELIN, Spectral Approximation of Linear Operators, Academic Press, New York, 1983.
- [15] M. DAUGE, Computations for Maxwell equations for the approximation of highly singular solutions, 2004, <http://perso.univ-rennes1.fr/monique.dauge/benchmax.html>.
- [16] E. B. DAVIES, Spectral enclosures and complex resonances for general self-adjoint operators, *LMS J. Comput. Math.*, 1 (1998), pp. 42–74.
- [17] ———, A hierarchical method for obtaining eigenvalue enclosures, *Math. Comp.*, 69 (2000), pp. 1435–1455.
- [18] E. B. DAVIES AND M. PLUM, Spectral pollution, *IMA J. Numer. Anal.*, 24 (2004), pp. 417–438.

- [19] A. ERN AND J.-L. GUERMOND, Theory and practice of finite elements, vol. 159 of Applied Mathematical Sciences, Springer-Verlag, New York, 2004.
- [20] V. GIRAULT AND P.-A. RAVIART, Finite element methods for Navier-Stokes equations, vol. 5 of Springer Series in Computational Mathematics, Springer-Verlag, Berlin, 1986. Theory and algorithms.
- [21] F. GOERISCH AND J. ALBRECHT, The convergence of a new method for calculating lower bounds to eigenvalues, in Equadiff 6 (Brno, 1985), vol. 1192 of Lecture Notes in Math., Springer, Berlin, 1986, pp. 303–308.
- [22] P. MONK, Finite element methods for Maxwell's equations, Clarendon Press, Cambridge, 2003.
- [23] G. STRANG AND G. FIX, An Analysis of the Finite Element Method, Prentice Hall, London, 1973.
- [24] S. ZIMMERMANN AND U. MERTINS, Variational bounds to eigenvalues of self-adjoint eigenvalue problems with arbitrary spectrum, *Z. Anal. Anwendungen*, 14 (1995), pp. 327–345.

DEPARTMENT OF MATHEMATICS AND STATISTICS, UNIVERSITY OF STRATHCLYDE, 26 RICHMOND STREET, GLASGOW G1 1XH, SCOTLAND

E-mail address: gabriel.barrenechea@strath.ac.uk

DEPARTMENT OF MATHEMATICS AND MAXWELL INSTITUTE FOR MATHEMATICAL SCIENCES, HERIOT-WATT UNIVERSITY, EDINBURGH, EH14 4AS, UK

E-mail address: L.Boulton@hw.ac.uk

DÉPARTEMENT DE MATHÉMATIQUES, UNIVERSITÉ DE FRANCHE-COMTÉ, BESANÇON, FRANCE

E-mail address: nboussai@univ-fcomte.fr