

## GENERALIZED SPECTRUM OF SECOND ORDER DIFFERENTIAL OPERATORS\*

TOMÁŠ GERGELITS<sup>†</sup>, BJØRN FREDRIK NIELSEN<sup>‡</sup>, AND ZDENĚK STRAKOŠ<sup>§</sup>

**Abstract.** We analyze the spectrum of the operator  $\Delta^{-1}[\nabla \cdot (K\nabla u)]$ , where  $\Delta$  denotes the Laplacian and  $K = K(x, y)$  is a symmetric tensor. Our main result shows that this spectrum can be derived from the spectral decomposition  $K = Q\Lambda Q^T$ , where  $Q = Q(x, y)$  is an orthogonal matrix and  $\Lambda = \Lambda(x, y)$  is a diagonal matrix. More precisely, provided that  $K$  is continuous, the spectrum equals the convex hull of the ranges of the diagonal function entries of  $\Lambda$ . The involved domain is assumed to be bounded and Lipschitz, and both homogeneous Dirichlet and homogeneous Neumann boundary conditions are considered. We study operators defined on infinite dimensional Sobolev spaces. Our theoretical investigations are illuminated by numerical experiments, using discretized problems. The results presented in this paper extend previous analyses which have addressed elliptic differential operators with scalar coefficient functions. Our investigation is motivated by both preconditioning issues (efficient numerical computations) and the need to further develop the spectral theory of second order PDEs (core analysis).

**Key words.** second order PDEs, generalized eigenvalues, spectrum, tensors, preconditioning

**AMS subject classifications.** 65F08, 65F15, 65N12, 35J99

**DOI.** 10.1137/20M1316159

**1. Introduction.** For simple domains, the eigenfunctions and eigenvalues of the Laplacian  $\Delta$  can be characterized in terms of trigonometric functions. Similar analytic information about the spectrum of general second order differential operators  $\nabla \cdot (K\nabla u)$  is not available. On the other hand, in [4, 10] the authors show that the generalized eigenvalue problem

$$\nabla \cdot (k\nabla u) = \lambda\Delta u \quad \text{for } (x, y) \in \Omega,$$

where  $k$  is a uniformly positive scalar function, can be analyzed in detail. More specifically, if  $k$  is continuous, then the range

$$k(\Omega) = \{k(x, y), (x, y) \in \Omega\}$$

of  $k$  is contained in the spectrum of the operator  $\Delta^{-1}[\nabla \cdot (k\nabla u)]$ . Furthermore, for discretized problems, assuming that  $k$  is bounded and piecewise continuous, the function values of  $k$  over the patches defined by the discretization basis functions provide accurate approximations of the generalized eigenvalues.

The main purpose of this paper is to extend the results published in [4, 10] to second order differential operators which involve a symmetric tensor, that is, to the

---

\*Received by the editors January 31, 2020; accepted for publication April 28, 2020; published electronically July 29, 2020.

<https://doi.org/10.1137/20M1316159>

**Funding:** The research of the first and third authors was supported by the Grant Agency of the Czech Republic under grant 17-04150J. The work of the second author was supported by the Research Council of Norway under project 239070.

<sup>†</sup>Institute of Computer Science of the Czech Academy of Sciences, 182 07 Prague 8, Czech Republic (gergelits@cs.cas.cz), and Faculty of Mathematics and Physics, Charles University, 186 75 Prague 8, Czech Republic (gergelits@karlin.mff.cuni.cz).

<sup>‡</sup>Faculty of Science and Technology, Norwegian University of Life Sciences, NO-1432 Ås, Norway (bjorn.f.nielsen@nmbu.no).

<sup>§</sup>Faculty of Mathematics and Physics, Charles University, 186 75 Prague 8, Czech Republic (strakos@karlin.mff.cuni.cz).

generalized eigenvalue problem

$$(1.1) \quad \begin{aligned} \nabla \cdot (K \nabla u) &= \lambda \Delta u && \text{for } (x, y) \in \Omega, \\ u &= 0 && \text{for } (x, y) \in \partial\Omega, \end{aligned}$$

where the open domain  $\Omega \subset \mathbb{R}^2$  is bounded and Lipschitz, and the real valued tensor function  $K : \Omega \rightarrow \mathbb{R}^{2 \times 2}$  is symmetric with its entries being bounded Lebesgue integrable functions and with the spectral decomposition

$$(1.2) \quad \begin{aligned} K(x, y) &= Q(x, y) \Lambda(x, y) Q^T(x, y), && (x, y) \in \Omega, \\ \Lambda(x, y) &= \begin{bmatrix} \kappa_1(x, y) & 0 \\ 0 & \kappa_2(x, y) \end{bmatrix}, && Q Q^T = Q^T Q = I. \end{aligned}$$

More specifically, defining the operators  $\mathcal{L}, \mathcal{A} : H_0^1(\Omega) \mapsto H^{-1}(\Omega)$  as

$$(1.3) \quad \langle \mathcal{L}\phi, \psi \rangle = \int_{\Omega} \nabla \phi \cdot \nabla \psi, \quad \phi, \psi \in H_0^1(\Omega),$$

$$(1.4) \quad \langle \mathcal{A}\phi, \psi \rangle = \int_{\Omega} K \nabla \phi \cdot \nabla \psi, \quad \phi, \psi \in H_0^1(\Omega),$$

we characterize the spectrum of the preconditioned operator

$$(1.5) \quad \mathcal{L}^{-1} \mathcal{A} : H_0^1(\Omega) \rightarrow H_0^1(\Omega),$$

defined as<sup>1</sup>

$$(1.6) \quad \text{sp}(\mathcal{L}^{-1} \mathcal{A}) \equiv \{ \lambda \in \mathbb{C}; \lambda \mathcal{L} - \mathcal{L}^{-1} \mathcal{A} \text{ does not have a bounded inverse} \}.$$

This paper proves the following result.

**THEOREM 1.1** (spectrum of the preconditioned operator). *Consider an open and bounded Lipschitz domain  $\Omega \subset \mathbb{R}^2$ . Assume that the tensor  $K$  is symmetric and continuous throughout the closure  $\bar{\Omega}$ . Then the spectrum of the operator  $\mathcal{L}^{-1} \mathcal{A}$ , defined in (1.3)–(1.6), equals*

$$(1.7) \quad \text{sp}(\mathcal{L}^{-1} \mathcal{A}) = \text{Conv}(\kappa_1(\bar{\Omega}) \cup \kappa_2(\bar{\Omega})),$$

where

$$(1.8) \quad \text{Conv}(\kappa_1(\bar{\Omega}) \cup \kappa_2(\bar{\Omega})) = \left[ \inf_{(x,y) \in \bar{\Omega}} \min_{i=1,2} \kappa_i(x, y), \sup_{(x,y) \in \bar{\Omega}} \max_{i=1,2} \kappa_i(x, y) \right].$$

Note that this theorem extends the results in [10] in several ways. It holds for second order differential operators with definite, indefinite, and semidefinite tensors. Moreover, instead of the inclusion proved for the scalar case in [10], it shows that the spectrum actually equals the interval (1.8) determined by  $K(x, y)$ .

Our theoretical study addresses operators defined on infinite dimensional Sobolev spaces. Numerical experiments suggest that even stronger properties, analogous to the scalar case analyzed in [4], hold for discretized problems.

<sup>1</sup>For operators defined on infinite dimensional Hilbert (Sobolev) spaces, the eigenvalues represent, in general, only a part of the spectrum. Therefore, the generalized eigenvalue problem (1.1) does not determine the whole spectrum (1.6).

Our theoretical results can be illustrated by the following experiment. We consider three test problems (1.1) with diagonal tensors (1.2) (i.e.,  $Q = I$ ) defined on the domain  $\Omega \equiv (0, 1) \times (0, 1)$ , where

$$(1.9) \quad \begin{aligned} \text{(P1)} : \quad & \kappa_1(x, y) = 1, & \kappa_2(x, y) &= 10, \\ \text{(P2)} : \quad & \kappa_1(x, y) = 1 + 0.5(x + y), & \kappa_2(x, y) &= 10 - 0.5(x + y), \\ \text{(P3)} : \quad & \kappa_1(x, y) = 1 + 3(x + y), & \kappa_2(x, y) &= 10 - 2(x + y) \end{aligned}$$

for  $(x, y) \in \Omega$ . We discretize the problem (1.1) using a uniform triangular mesh with piecewise linear discretization basis functions; see [4] for the scalar case analogy. Figure 1 presents the eigenvalues of the resulting discrete generalized eigenvalue problem of size 381. We observe that the spectrum of the discretized problem covers not only the union of the ranges  $\kappa_1(\bar{\Omega}) \cup \kappa_2(\bar{\Omega})$ , but in the case that  $\kappa_1(\bar{\Omega})$  and  $\kappa_2(\bar{\Omega})$  do not overlap, it surprisingly covers the whole interval (1.8).

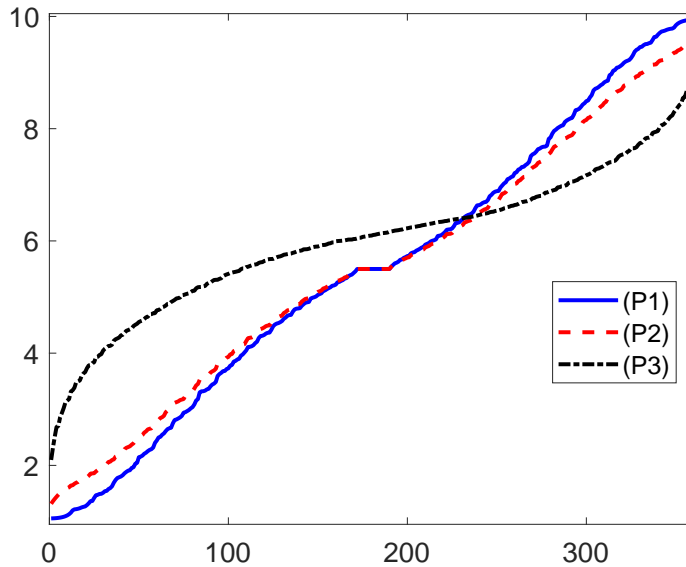


FIG. 1. Eigenvalues of the discretized problems (P1)–(P3), defined in (1.9), spread over the entire interval  $[1, 10]$ , while the ranges of entries of the diagonal tensor are the following: (P1):  $\kappa_1(\bar{\Omega}) = 1, \kappa_2(\bar{\Omega}) = 10$ ; (P2):  $\kappa_1(\bar{\Omega}) = [1, 2], \kappa_2(\bar{\Omega}) = [9, 10]$ ; (P3):  $\kappa_1(\bar{\Omega}) = [1, 7], \kappa_2(\bar{\Omega}) = [6, 10]$ . Horizontal axis: the indices of the increasingly ordered eigenvalues. Vertical axis: the size of the eigenvalues.

Since  $\langle \mathcal{A}u, v \rangle = \langle \mathcal{A}v, u \rangle$  for all  $u, v \in H_0^1(\Omega)$ , which is a consequence of the symmetry of the tensor  $K$ , the preconditioned operator (1.5) is self-adjoint with respect to the inner product associated with the Laplacian:

$$(1.10) \quad (u, v)_{\mathcal{L}} \equiv \langle \mathcal{L}u, v \rangle = \int_{\Omega} \nabla u \cdot \nabla v, \quad u, v \in H_0^1(\Omega),$$

$$(1.11) \quad (\mathcal{L}^{-1}\mathcal{A}u, v)_{\mathcal{L}} = \langle \mathcal{A}u, v \rangle = \langle \mathcal{A}v, u \rangle = (\mathcal{L}^{-1}\mathcal{A}v, u)_{\mathcal{L}}.$$

Consequently,  $\text{sp}(\mathcal{L}^{-1}\mathcal{A}) \subset \mathbb{R}$ . The inner product (1.10) defines the norm

$$\|u\|_{\mathcal{L}}^2 \equiv (u, u)_{\mathcal{L}} = \langle \mathcal{L}u, u \rangle = \int_{\Omega} \|\nabla u\|^2 = \int_{\Omega} \|u_x\|_2^2 + \|u_y\|_2^2, \quad u \in H_0^1(\Omega),$$

used in the proofs below.

The convergence behavior of the conjugate gradient (CG) method is determined by the spectral distribution functions of the involved linear systems; see, e.g., [5, 8]. Hence, the analysis presented in this paper can be employed to better understand the performance of CG when the inverse of the Laplacian (or some variant incorporating it) is applied as the preconditioner to solve discretized second order elliptic PDEs; see [4] for a discussion of this topic. Also, constant-coefficient preconditioners may be of particular interest when the isogeometric analysis (IgA) approach is employed to discretize both PDEs and the involved computational domains in terms of B-splines [6, 7, 12].

This paper is organized in the following way. For clarity of exposition, we restrict ourselves in sections 2 and 3 to problems with diagonal tensors. In section 2 we present auxiliary lemmas generalizing, step by step, the results in [10]. Section 3 contains the proof of the main result for problems with diagonal tensors, and in section 4 we generalize the lemmas from previous sections to nondiagonal symmetric tensors and give the proof of the main result, Theorem 1.1. In section 5 we comment on problems with homogeneous Neumann boundary conditions. The numerical experiments in section 6 illustrate the results of the analysis, and the text closes with a brief discussion of some open problems in section 7.

**2. Auxiliary results.** We will start with considering diagonal tensors, i.e.,

$$(2.1) \quad K(x, y) = \begin{bmatrix} \kappa_1(x, y) & 0 \\ 0 & \kappa_2(x, y) \end{bmatrix}.$$

This will allow us to explain with full clarity the main difference between the scalar case studied in [4, 10] and the tensor case analyzed in this paper.

**2.1. Function values at points of continuity belong to the spectrum.**

The following lemma generalizes statement (a) in Theorem 3.1 in [10].

**LEMMA 2.1.** *Assume that  $K$  is a diagonal tensor, where the entries  $\kappa_1$  and  $\kappa_2$  are bounded and Lebesgue integrable functions on  $\Omega$ . The following holds for  $i = 1, 2$ : If  $\kappa_i$  is continuous at  $(x_0, y_0) \in \Omega$ , then*

$$\kappa_i(x_0, y_0) \in \text{sp}(\mathcal{L}^{-1}\mathcal{A}).$$

*Proof.* Assume that  $\kappa_1$  is continuous at  $(x_0, y_0)$ , and let

$$\lambda \equiv \kappa_1(x_0, y_0).$$

We will construct parametrized functions  $v_r$  and  $u_r = (\lambda\mathcal{I} - \mathcal{L}^{-1}\mathcal{A})v_r$  such that

$$(2.2) \quad \lim_{r \rightarrow 0} \|v_r\|_{\mathcal{L}} \neq 0 \quad \text{and} \quad \lim_{r \rightarrow 0} \|u_r\|_{\mathcal{L}} = 0,$$

which is not possible if  $\lambda\mathcal{I} - \mathcal{L}^{-1}\mathcal{A}$  has a bounded inverse:  $v_r = (\lambda\mathcal{I} - \mathcal{L}^{-1}\mathcal{A})^{-1}u_r$  and  $\lim_{r \rightarrow 0} \|u_r\|_{\mathcal{L}} = 0$  imply that  $\lim_{r \rightarrow 0} \|v_r\|_{\mathcal{L}} = 0$ . (The norm  $\|\cdot\|_{\mathcal{L}}$  is the norm induced by the inner product (1.10).)

The functions  $v_r$  can be constructed, e.g., in the following way. Consider, for a sufficiently small  $r > 0$ , the following closed neighborhood of the point  $(x_0, y_0)$ :

$$(2.3) \quad R_r = [x_0 - r^2, x_0 + r^2] \times [y_0 - r, y_0 + r] \subset \Omega.$$

For  $(x, y) \in R_r$ , define

$$(2.4) \quad v_r(x, y) = \sqrt{r} \min \left\{ 1 - \frac{|x-x_0|}{r^2}, \frac{1}{r} - \frac{|y-y_0|}{r^2} \right\},$$

and  $v_r(x, y) = 0$  otherwise. It can be verified that (see Appendix A)

$$(2.5) \quad \begin{aligned} 4 - 4r &\leq \|(v_r)_x\|_{L^2(\Omega)}^2 \leq 4, \\ &\|(v_r)_y\|_{L^2(\Omega)}^2 \leq 4r. \end{aligned}$$

Consequently,

$$(2.6) \quad \lim_{r \rightarrow 0} \|v_r\|_{\mathcal{L}} = \lim_{r \rightarrow 0} \left( \|(v_r)_x\|_{L^2(\Omega)}^2 + \|(v_r)_y\|_{L^2(\Omega)}^2 \right)^{1/2} = 2.$$

Considering

$$(2.7) \quad u_r = (\lambda \mathcal{I} - \mathcal{L}^{-1} \mathcal{A})v_r, \quad \text{i.e.,} \quad \mathcal{L}u_r = (\lambda \mathcal{L} - \mathcal{A})v_r,$$

we get

$$\begin{aligned} \|u_r\|_{\mathcal{L}}^2 &= \langle \mathcal{L}u_r, u_r \rangle = \langle (\lambda \mathcal{L} - \mathcal{A})v_r, u_r \rangle \\ &= \int_{\Omega} (\lambda I - K) \nabla v_r \cdot \nabla u_r \\ &\leq \left( \int_{\Omega} |(\lambda I - K) \nabla v_r|^2 \right)^{1/2} \|u_r\|_{\mathcal{L}}. \end{aligned}$$

Using that  $\text{supp}(v_r) = R_r$  and (2.5),

$$\begin{aligned} \|u_r\|_{\mathcal{L}}^2 &\leq \|(\lambda - \kappa_1)(v_r)_x\|_{L^2(\Omega)}^2 + \|(\lambda - \kappa_2)(v_r)_y\|_{L^2(\Omega)}^2 \\ &\leq 4 \sup_{(x,y) \in R_r} |\kappa_1(x_0, y_0) - \kappa_1(x, y)|^2 + 4r(\|\kappa_1\|_{L^\infty(\Omega)} + \|\kappa_2\|_{L^\infty(\Omega)})^2, \end{aligned}$$

and from the continuity of  $\kappa_1(x, y)$  at  $(x_0, y_0)$ ,

$$(2.8) \quad \lim_{r \rightarrow 0} \|u_r\|_{\mathcal{L}} = 0.$$

From (2.6) and (2.8) we conclude that we can construct functions  $v_r$  and  $u_r = (\lambda \mathcal{I} - \mathcal{L}^{-1} \mathcal{A})v_r$  such that (2.2) holds. We conclude that  $\kappa_1(x_0, y_0) \mathcal{I} - \mathcal{L}^{-1} \mathcal{A}$  cannot have a bounded inverse.

The proof that  $\kappa_2(x_0, y_0)$  belongs to the spectrum if  $\kappa_2$  is continuous at  $(x_0, y_0)$  is trivially analogous.  $\square$

If  $\kappa_i \in \mathcal{C}(\Omega)$ ,  $i = 1, 2$ , then Lemma 2.1 gives a diagonal-tensor-case analogy of Theorem 3.1, statement (b), in [10]. As is shown next, in the tensor case the spectrum of the preconditioned operator  $\mathcal{L}^{-1} \mathcal{A}$  can, however, also contain numbers that do not belong to any of the individual ranges of the functions  $\kappa_1$  and  $\kappa_2$ .

**2.2. Disjoint ranges extend the spectrum.** An unexpected case occurs when the ranges of  $\kappa_1$  and  $\kappa_2$  are disjoint:

$$\kappa_1(\overline{\Omega}) \cap \kappa_2(\overline{\Omega}) = \emptyset.$$

We begin by presenting the following facts that will be used in the proofs.

**2.2.1. Dirichlet problem for the wave equation.** Note that for any integer  $n$ ,

$$(2.9) \quad \phi(x, y) = \sin(n\pi cl^{-1}(y - y_0)) \sin(n\pi l^{-1}(x - x_0))$$

solves the following Dirichlet problem for the wave equation:

$$(2.10) \quad \begin{aligned} \phi_{yy} &= c^2 \phi_{xx} && \text{in } \Sigma_l, \\ \phi &= 0 && \text{on } \partial\Sigma_l, \end{aligned}$$

where  $l$  is a positive constant which determines the size of the solution domain

$$\Sigma_l = (x_0, x_0 + l) \times (y_0, y_0 + l/c),$$

and  $c > 0$  is arbitrary. We conclude that this Dirichlet problem has infinitely many nontrivial solutions. It is also clear that  $\Sigma_l$  can be made as small as needed by choosing  $l > 0$  sufficiently small.

**2.2.2. Tensors constant on an open subdomain.** Consider the generalized eigenvalue problem (1.1) with a diagonal tensor  $K(x, y)$  (2.1) that is constant on an open subdomain  $S \subset \Omega$ . Then we get the following lemma.

LEMMA 2.2. *Consider a diagonal tensor (2.1), where the bounded and Lebesgue integrable functions  $\kappa_i$ ,  $i = 1, 2$ , are constant on an open subdomain  $S \subset \Omega$ . Assuming that*

$$(2.11) \quad \sup_{(x,y) \in \Omega} \kappa_1(x, y) < \inf_{(x,y) \in \Omega} \kappa_2(x, y),$$

the following closed interval belongs to the spectrum of  $\mathcal{L}^{-1}\mathcal{A}$ :

$$(2.12) \quad \left[ \sup_{(x,y) \in \Omega} \kappa_1(x, y), \inf_{(x,y) \in \Omega} \kappa_2(x, y) \right] \subset \text{sp}(\mathcal{L}^{-1}\mathcal{A}).$$

The analogous statement obviously holds with interchanging the roles of  $\kappa_1$  and  $\kappa_2$  in (2.11) and (2.12).

*Proof.* Consider an arbitrary fixed point  $(x_0, y_0) \in S$ . For any fixed  $c > 0$ , there exists  $l > 0$  such that

$$\Sigma_l \equiv (x_0, x_0 + l) \times (y_0, y_0 + l/c) \subset S.$$

Since  $K(x, y)$  is constant on  $\Sigma_l$ , we can rewrite (1.1) as

$$(2.13) \quad (\lambda - \bar{k}_1)v_{xx} + (\lambda - \bar{k}_2)v_{yy} = 0 \quad \text{in } \Sigma_l,$$

where  $\bar{k}_1$  and  $\bar{k}_2$  are constants and

$$K(x, y) = \begin{bmatrix} \bar{k}_1 & 0 \\ 0 & \bar{k}_2 \end{bmatrix}, \quad (x, y) \in \Sigma_l.$$

Consider an arbitrary  $\lambda$  in the interval  $(\bar{k}_1, \bar{k}_2)$ . Then (2.13) represents, with

$$c^2 = \frac{\lambda - \bar{k}_1}{\bar{k}_2 - \lambda} > 0,$$

the wave equation (2.10). Taking any nontrivial solution  $\phi$  of (2.10), the function  $v$  defined on  $\Omega$  as

$$v(x, y) = \begin{cases} \phi(x, y), & (x, y) \in \Sigma_l, \\ 0, & (x, y) \notin \Sigma_l, \end{cases}$$

solves the weak form of the generalized eigenvalue problem (1.1). We conclude that  $(\bar{k}_1, \bar{k}_2) \subset \text{sp}(\mathcal{L}^{-1}\mathcal{A})$ .

Since, by construction,

$$(2.14) \quad \bar{k}_1 \leq \sup_{(x,y) \in \Omega} \kappa_1(x, y) < \inf_{(x,y) \in \Omega} \kappa_2(x, y) \leq \bar{k}_2,$$

it remains to prove that if the equality is attained at any side of (2.14), then the associated  $\bar{k}_i$ ,  $i = 1$  and/or  $i = 2$ , also belongs to the spectrum of  $\mathcal{L}^{-1}\mathcal{A}$ . But this is trivially true using Lemma 2.1 because  $\bar{k}_i$  is a function value of  $\kappa_i(x, y)$  at  $\Sigma_l$ , where  $\kappa_i$  is constant and therefore continuous.  $\square$

Lemma 2.2 shows that, under the given assumptions, the whole closed interval determined by the extremal points of the ranges of  $\kappa_1$  and  $\kappa_2$  belong to the spectrum of  $\mathcal{L}^{-1}\mathcal{A}$ . Consequently, when the ranges of  $\kappa_1$  and  $\kappa_2$  are disjoint, the spectrum of  $\mathcal{L}^{-1}\mathcal{A}$  contains also the interval between them. Please note that here it is not assumed that  $K$  is continuous throughout the closure  $\bar{\Omega}$  and that the subdomain  $S$  is of an arbitrarily small size.

**2.2.3. Tensors continuous at least at a single point.** The following lemma refines further the assumptions under which the statement of Lemma 2.2 holds.

LEMMA 2.3. *Assume that the diagonal tensor (2.1) with the bounded and Lebesgue integrable functions  $\kappa_i$ ,  $i = 1, 2$ , is continuous (at least) at a single point in  $\Omega$ . If*

$$(2.15) \quad \sup_{(x,y) \in \Omega} \kappa_1(x, y) < \inf_{(x,y) \in \Omega} \kappa_2(x, y),$$

then the following closed interval belongs to the spectrum of  $\mathcal{L}^{-1}\mathcal{A}$ :

$$(2.16) \quad \left[ \sup_{(x,y) \in \Omega} \kappa_1(x, y), \inf_{(x,y) \in \Omega} \kappa_2(x, y) \right] \subset \text{sp}(\mathcal{L}^{-1}\mathcal{A}).$$

The analogous statement obviously holds with interchanging the roles of  $\kappa_1$  and  $\kappa_2$  in (2.15) and (2.16).

*Proof.* We will prove the statement by contradiction. Consider

$$\lambda \in \left[ \sup_{(x,y) \in \Omega} \kappa_1(x, y), \inf_{(x,y) \in \Omega} \kappa_2(x, y) \right]$$

such that  $\lambda \notin \text{sp}(\mathcal{L}^{-1}\mathcal{A})$ , i.e., such that the operator  $\mathcal{L}^{-1}\mathcal{A} - \lambda\mathcal{I}$  has a bounded inverse.

Let  $(x_0, y_0) \in \Omega$  be the point of continuity of the tensor  $K(x, y)$ . Applying Lemma 2.2 to the preconditioned operator  $\mathcal{L}^{-1}\mathcal{A}_l$ , where  $\mathcal{A}_l$  is defined for any sufficiently small  $l$  by

$$\langle \mathcal{A}_l \phi, \psi \rangle \equiv \int_{\Omega} K_l \nabla \phi \cdot \nabla \psi, \quad \phi, \psi \in H_0^1(\Omega),$$

and  $K_l(x, y)$  is a local modification of  $K$ ,

$$K_l(x, y) \equiv \begin{cases} K(x_0, y_0), & (x, y) \in S_l, \\ K(x, y), & (x, y) \in \Omega \setminus S_l, \end{cases}$$

$$S_l = (x_0, x_0 + l) \times (y_0, y_0 + l),$$

yields that

$$(2.17) \quad \lambda \in \text{sp}(\mathcal{L}^{-1}\mathcal{A}_l).$$

On the other hand, since we assume that  $\mathcal{L}^{-1}\mathcal{A} - \lambda\mathcal{I}$  is invertible,

$$\begin{aligned} \mathcal{L}^{-1}\mathcal{A}_l - \lambda\mathcal{I} &= (\mathcal{L}^{-1}\mathcal{A} - \lambda\mathcal{I}) + (\mathcal{L}^{-1}\mathcal{A}_l - \mathcal{L}^{-1}\mathcal{A}) \\ &= (\mathcal{L}^{-1}\mathcal{A} - \lambda\mathcal{I})[\mathcal{I} + (\mathcal{L}^{-1}\mathcal{A} - \lambda\mathcal{I})^{-1}\mathcal{L}^{-1}(\mathcal{A}_l - \mathcal{A})]. \end{aligned}$$

In Appendix B we prove that for sufficiently small  $l > 0$

$$(2.18) \quad \|(\mathcal{L}^{-1}\mathcal{A} - \lambda\mathcal{I})^{-1}\mathcal{L}^{-1}(\mathcal{A}_l - \mathcal{A})\|_{\mathcal{L}} < 1, \quad \square$$

and the Neumann series argument therefore ensures that  $\mathcal{L}^{-1}\mathcal{A}_l - \lambda\mathcal{I}$  has a bounded inverse. Consequently,  $\lambda \notin \text{sp}(\mathcal{L}^{-1}\mathcal{A}_l)$ , which contradicts (2.17). (Inequality (2.18) holds due to the assumption that  $\lambda \notin \text{sp}(\mathcal{L}^{-1}\mathcal{A})$  and due to the continuity of  $K(x, y)$  at the point  $(x_0, y_0)$ . See Appendix B for further details.)

It is worth noting that the statement of Lemma 2.3 requires continuity of the tensor  $K$  only at an arbitrary *single point* belonging to  $\Omega$ .

**3. Continuous diagonal tensors.** We first complement Lemma 2.1, and Theorem 3.1 in [10], by proving the “reverse inclusion.”

### 3.1. The spectrum is a subset of the extremal interval.

LEMMA 3.1. *Assume that the diagonal tensor (2.1) is continuous throughout the closure  $\bar{\Omega}$ . Then*

$$(3.1) \quad \text{sp}(\mathcal{L}^{-1}\mathcal{A}) \subset \text{Conv}(\kappa_1(\bar{\Omega}) \cup \kappa_2(\bar{\Omega})).$$

*Proof.* Using the self-adjointness (1.11) of the operator  $\mathcal{L}^{-1}\mathcal{A}$ , we can take the standard results from the theory of self-adjoint operators (see, e.g., [3, section 6.5]) and conclude that the spectrum of  $\mathcal{L}^{-1}\mathcal{A}$  is real and that

$$(3.2) \quad \begin{aligned} \text{sp}(\mathcal{L}^{-1}\mathcal{A}) &\subset \left[ \inf_{u \in H_0^1(\Omega)} \frac{(\mathcal{L}^{-1}\mathcal{A}u, u)_{\mathcal{L}}}{(u, u)_{\mathcal{L}}}, \sup_{u \in H_0^1(\Omega)} \frac{(\mathcal{L}^{-1}\mathcal{A}u, u)_{\mathcal{L}}}{(u, u)_{\mathcal{L}}} \right] \\ &= \left[ \inf_{u \in H_0^1(\Omega)} \frac{\langle \mathcal{A}u, u \rangle}{\langle \mathcal{L}u, u \rangle}, \sup_{u \in H_0^1(\Omega)} \frac{\langle \mathcal{A}u, u \rangle}{\langle \mathcal{L}u, u \rangle} \right]. \end{aligned}$$

Moreover, the endpoints of this interval are contained in the spectrum.

It remains to bound

$$(3.3) \quad \frac{\langle \mathcal{A}u, u \rangle}{\langle \mathcal{L}u, u \rangle}$$

in terms of the extreme values of the scalar functions  $\kappa_1$  and  $\kappa_2$ . Since  $u_x^2(x, y) \geq 0$  and  $u_y^2(x, y) \geq 0$ , we can bound (3.3) as follows:

$$\begin{aligned}
 \sup_{u \in H_0^1(\Omega)} \frac{\langle \mathcal{A}u, u \rangle}{\langle \mathcal{L}u, u \rangle} &= \sup_{u \in H_0^1(\Omega)} \frac{\int_{\Omega} K \nabla u \cdot \nabla u}{\int_{\Omega} \|\nabla u\|^2} = \sup_{u \in H_0^1(\Omega)} \frac{\int_{\Omega} \kappa_1 u_x^2 + \kappa_2 u_y^2}{\int_{\Omega} \|\nabla u\|^2} \\
 &\leq \sup_{u \in H_0^1(\Omega)} \frac{\int_{\Omega} \sup_{(x,y) \in \Omega} \max_{i=1,2} \{\kappa_i(x, y)\} \|\nabla u\|^2}{\int_{\Omega} \|\nabla u\|^2} \\
 (3.4) \quad &\leq \sup_{(x,y) \in \Omega} \max_{i=1,2} \{\kappa_i(x, y)\}.
 \end{aligned}$$

Similarly,

$$\inf_{u \in H_0^1(\Omega)} \frac{\langle \mathcal{A}u, u \rangle}{\langle \mathcal{L}u, u \rangle} \geq \inf_{(x,y) \in \Omega} \min_{i=1,2} \{\kappa_i(x, y)\}.$$

For  $K(x, y)$  continuous on  $\bar{\Omega}$ , the infimum and supremum of its components  $\kappa_1(x, y)$  and  $\kappa_2(x, y)$  are attained. Please notice that there is no assumption about the positive (negative) definiteness of  $K$ .  $\square$

We are now ready to prove Theorem 1.1 for continuous diagonal tensors.

**3.2. Main result — diagonal tensors.**

**THEOREM 3.2.** *Consider an open and bounded Lipschitz domain  $\Omega \subset \mathbb{R}^2$ . If the diagonal tensor (2.1) is continuous throughout the closure  $\bar{\Omega}$ , then*

$$\text{sp}(\mathcal{L}^{-1}\mathcal{A}) = \text{Conv}(\kappa_1(\bar{\Omega}) \cup \kappa_2(\bar{\Omega})).$$

*Proof.* Assume that the diagonal tensor  $K(x, y)$  is continuous throughout  $\bar{\Omega}$ . Then, by Lemmas 2.1 and 2.3,

$$\text{Conv}(\kappa_1(\Omega) \cup \kappa_2(\Omega)) \subset \text{sp}(\mathcal{L}^{-1}\mathcal{A}),$$

and due to the continuity of  $K(x, y)$  and the fact that  $\text{sp}(\mathcal{L}^{-1}\mathcal{A})$  is a closed set (see, e.g., [11]),

$$\text{Conv}(\kappa_1(\bar{\Omega}) \cup \kappa_2(\bar{\Omega})) \subset \text{sp}(\mathcal{L}^{-1}\mathcal{A}).$$

Finally, by Lemma 3.1,

$$\text{sp}(\mathcal{L}^{-1}\mathcal{A}) \subset \text{Conv}(\kappa_1(\bar{\Omega}) \cup \kappa_2(\bar{\Omega})),$$

which gives the statement.  $\square$

**4. Proof of Theorem 1.1.** It remains to revisit and complete the arguments given above for the general self-adjoint operator in (1.1). Consider the general symmetric tensor

$$(4.1) \quad K(x, y) = \begin{bmatrix} k_1(x, y) & k_3(x, y) \\ k_3(x, y) & k_2(x, y) \end{bmatrix},$$

where  $k_1, k_2$ , and  $k_3$  are bounded and Lebesgue integrable functions defined on  $\Omega$ , with the spectral decomposition

$$(4.2) \quad K(x, y) = Q(x, y) \begin{bmatrix} \kappa_1(x, y) & 0 \\ 0 & \kappa_2(x, y) \end{bmatrix} Q^T(x, y);$$

Downloaded 04/30/24 to 195.113.30.215 . Redistribution subject to SIAM license or copyright; see https://pubs.siam.org/terms-privacy

see (1.2).

The structure of the proof of Theorem 1.1 is fully analogous to the proof of Theorem 3.2 formulated for diagonal tensors. We will now restate the associated lemmas for the general case and comment on the technical differences that must be considered.

For convenience, we will use, when appropriate, the column vector notation

$$\mathbf{w} = (x, y)^T, \quad (x, y) \in \Omega,$$

and for any function  $f$  defined on  $\Omega$  its gradient  $\nabla f$  will be considered as a column vector.

LEMMA 4.1 (see Lemma 2.1). *Consider the symmetric tensor (4.1) with the spectral decomposition (4.2). If the tensor  $K$  is continuous at  $(x_0, y_0) \in \Omega$ , then*

$$\kappa_i(x_0, y_0) \in \text{sp}(\mathcal{L}^{-1}\mathcal{A}), \quad i = 1, 2.$$

*Proof.* We will use the following notation for the spectral decomposition of  $K(x, y)$  at the point of continuity  $(x_0, y_0)$ :

$$\begin{aligned} K_0 &\equiv K(x_0, y_0) = Q_0 \Lambda_0 Q_0^T, & Q_0 &\equiv Q(x_0, y_0), & Q_0^T Q_0 &= I, \\ \Lambda_0 &\equiv \Lambda(x_0, y_0) = \text{diag}(\kappa_1(x_0, y_0), \kappa_2(x_0, y_0)). \end{aligned}$$

Simple algebraic computations give that, for any  $(x, y) \in \Omega$ ,

$$(4.3) \quad \kappa_1 = \frac{1}{2}(k_1 + k_2 + \sqrt{D}), \quad \kappa_2 = \frac{1}{2}(k_1 + k_2 - \sqrt{D}),$$

where  $D = (k_1 - k_2)^2 + 4k_3^2$ . Therefore, at any point of continuity of the tensor  $K(x, y)$ , the functions  $\kappa_1(x, y)$  and  $\kappa_2(x, y)$  are also continuous.

For sufficiently small  $r$ , consider the closed neighborhood  $R_r$  defined in (2.3) and its counterpart defined as

$$S_r = \{Q_0 \mathbf{z} \mid \mathbf{z} \in R_r\},$$

where the choice of  $r$  in (2.3) ensures that both  $R_r \subset \Omega$  and  $S_r \subset \Omega$ . Consider the functions

$$\tilde{v}_r(\mathbf{w}) \equiv v_r(Q_0^T \mathbf{w}), \quad \mathbf{w} \in \Omega,$$

where  $v_r$  is defined in (2.4). Since  $|\det Q| = 1$ , the change of variables gives

$$(4.4) \quad \|\tilde{v}_r\|_{\mathcal{L}}^2 = \int_{S_r} \|\nabla \tilde{v}_r(\mathbf{w})\|^2 d\mathbf{w} = \int_{R_r} \|\nabla v_r(\mathbf{z})\|^2 d\mathbf{z} = \|v_r\|_{\mathcal{L}}^2,$$

and, from (2.6),

$$(4.5) \quad \lim_{r \rightarrow 0} \|\tilde{v}_r\|_{\mathcal{L}} = 2 \neq 0.$$

Analogously to (2.7) we consider

$$u_r \equiv (\lambda \mathcal{I} - \mathcal{L}^{-1} \mathcal{A}) \tilde{v}_r, \quad \lambda \equiv \kappa_1(x_0, y_0),$$

with the norm

$$(4.6) \quad \|u_r\|_{\mathcal{L}}^2 = \int_{\Omega} (\lambda I - K) \nabla \tilde{v}_r \cdot \nabla u_r$$

$$(4.7) \quad = \int_{S_r} (\lambda I - K_0) \nabla \tilde{v}_r \cdot \nabla u_r + \int_{S_r} (K_0 - K) \nabla \tilde{v}_r \cdot \nabla u_r.$$

Our goal is to show that if  $\lambda \notin \text{sp}(\mathcal{L}^{-1}\mathcal{A})$ , then  $\lim_{r \rightarrow 0} \|u_r\|_{\mathcal{L}} = 0$ , which contradicts (4.5). Concerning the second integral in (4.7),

$$\int_{S_r} (K_0 - K) \nabla \tilde{v}_r \cdot \nabla u_r \leq \sup_{\mathbf{w} \in S_r} \|K_0 - K(\mathbf{w})\| \|\tilde{v}_r\|_{\mathcal{L}} \|u_r\|_{\mathcal{L}}.$$

Using the continuity of  $K(x, y)$  at the point  $(x_0, y_0)$  and the fact that  $\|\tilde{v}_r\|_{\mathcal{L}} \|u_r\|_{\mathcal{L}}$  is bounded, the second integral on the right-hand side of (4.7) vanishes as  $r \rightarrow 0$ . For the remaining term in (4.7), we find that

$$\begin{aligned} \int_{S_r} (\lambda I - K_0) \nabla \tilde{v}_r \cdot \nabla u_r &= \int_{S_r} Q_0 (\lambda I - \Lambda_0) Q_0^T \nabla \tilde{v}_r \cdot \nabla u_r \\ &\leq \left( \int_{S_r} \|Q_0 (\lambda I - \Lambda_0) Q_0^T \nabla \tilde{v}_r\|^2 \right)^{1/2} \|u_r\|_{\mathcal{L}}. \end{aligned}$$

Applying the chain rule gives  $\nabla \tilde{v}_r(\mathbf{w}) = Q_0 \nabla v_r(Q_0^T \mathbf{w}) = Q_0 \nabla v_r(\mathbf{z})$ , which together with the orthogonality of  $Q_0$  gives (considering  $\lambda = \kappa_1(x_0, y_0)$ )

$$\begin{aligned} \int_{S_r} \|Q_0 (\lambda I - \Lambda_0) Q_0^T \nabla \tilde{v}_r\|^2 &= \int_{S_r} \|(\lambda I - \Lambda_0) \nabla v_r(Q_0^T \mathbf{w})\|^2 \\ &= \int_{R_r} \|(\lambda I - \Lambda_0) \nabla v_r(\mathbf{z})\|^2 \\ &= \int_{R_r} \|(\lambda - \kappa_2(x_0, y_0)) (v_r)_y(\mathbf{z})\|^2 \\ &\leq |\lambda - \kappa_2(x_0, y_0)| \|(v_r)_y\|_{L^2(\Omega)}^2, \end{aligned}$$

where the upper bound vanishes as  $r \rightarrow 0$  due to (2.5).

The proof that  $\kappa_2(x_0, y_0)$  belongs to the spectrum of the preconditioned operator, provided that the assumptions of the lemma hold, is trivially analogous.  $\square$

The remaining part of the proof of Theorem 1.1 is a straightforward extension of the analysis presented in section 3.

LEMMA 4.2 (see Lemma 2.2). *Consider a symmetric tensor (4.1) with bounded and Lebesgue integrable functions  $k_i$ ,  $i = 1, 2, 3$ , which are constant on an open subdomain  $S \subset \Omega$ . Assuming that*

$$(4.8) \quad \sup_{(x,y) \in \Omega} \kappa_1(x, y) < \inf_{(x,y) \in \Omega} \kappa_2(x, y),$$

*the following closed interval belongs to the spectrum of  $\mathcal{L}^{-1}\mathcal{A}$ :*

$$(4.9) \quad \left[ \sup_{(x,y) \in \Omega} \kappa_1(x, y), \inf_{(x,y) \in \Omega} \kappa_2(x, y) \right] \subset \text{sp}(\mathcal{L}^{-1}\mathcal{A}).$$

*The analogous statement obviously holds with interchanging the roles of  $\kappa_1$  and  $\kappa_2$  in (4.8) and (4.9).*

*Proof.* Since  $K(x, y)$  and its spectral decomposition  $K = \bar{Q}\bar{\Lambda}\bar{Q}^T$  are constant on  $S$ , the change of variables  $\mathbf{w} = \bar{Q}\mathbf{z}$  transforms the eigenvalue problem (1.1) in the subdomain  $S$  to the form

$$\nabla_{\mathbf{z}} \cdot (\bar{\Lambda}\nabla_{\mathbf{z}}v) = \lambda\Delta_{\mathbf{z}}v \quad \text{in } R = \{\bar{Q}^T\mathbf{w} \mid \mathbf{w} \in S\},$$

where the diagonal tensor  $\bar{\Lambda}$  is constant. Employing the argument used to prove Lemma 2.2 finishes the proof.  $\square$

LEMMA 4.3 (see Lemma 2.3). *Assume that the symmetric tensor (4.1) with the bounded and Lebesgue integrable functions  $\kappa_i$ ,  $i = 1, 2, 3$ , is continuous at least at a single point in  $\Omega$ . Assuming that*

$$(4.10) \quad \sup_{(x,y) \in \Omega} \kappa_1(x, y) < \inf_{(x,y) \in \Omega} \kappa_2(x, y),$$

*the following closed interval belongs to the spectrum of  $\mathcal{L}^{-1}\mathcal{A}$ :*

$$(4.11) \quad \left[ \sup_{(x,y) \in \Omega} \kappa_1(x, y), \inf_{(x,y) \in \Omega} \kappa_2(x, y) \right] \subset \text{sp}(\mathcal{L}^{-1}\mathcal{A}).$$

*The analogous statement obviously holds with interchanging the roles of  $\kappa_1$  and  $\kappa_2$  in (4.10) and (4.11).*

*Proof.* The proof is fully analogous to the proof of Lemma 2.3.  $\square$

LEMMA 4.4 (see Lemma 3.1). *Let the symmetric tensor (4.1) be continuous throughout the closure  $\bar{\Omega}$ . Then*

$$\text{sp}(\mathcal{L}^{-1}\mathcal{A}) \subset \text{Conv}(\kappa_1(\bar{\Omega}) \cup \kappa_2(\bar{\Omega})).$$

*Proof.* The proof is analogous to the proof of Lemma 3.1 with the argument used in the derivation of (3.4) now written in the form

$$K\nabla u \cdot \nabla u = \Lambda Q^T \nabla u \cdot Q^T \nabla u \leq \sup_{\mathbf{w} \in \Omega} \max_{i=1,2} \{\kappa_i(\mathbf{w})\} \|Q^T \nabla u\|^2.$$

Due to the orthogonality of  $Q$  we get

$$\int_{\Omega} K\nabla u \cdot \nabla u \leq \sup_{\mathbf{w} \in \Omega} \max_{i=1,2} \{\kappa_i(\mathbf{w})\} \int_{\Omega} \|\nabla u\|^2,$$

and, similarly,

$$\inf_{\mathbf{w} \in \Omega} \max_{i=1,2} \{\kappa_i(\mathbf{w})\} \int_{\Omega} \|\nabla u\|^2 \leq \int_{\Omega} K\nabla u \cdot \nabla u. \quad \square$$

The proof of Theorem 1.1 is completed by the combination of Lemmas 4.1 to 4.4; see the proof of Theorem 3.2.

**5. Neumann boundary conditions.** Theorem 1.1 also holds for generalized eigenvalue problems with homogeneous Neumann boundary conditions:

$$(5.1) \quad \begin{aligned} \nabla \cdot (K\nabla u) &= \lambda\Delta u & \text{for } (x, y) \in \Omega, \\ \nabla u \cdot \mathbf{n} &= 0 & \text{for } (x, y) \in \partial\Omega, \end{aligned}$$

where  $\mathbf{n}$  denotes the outwards pointing unit normal vector of  $\partial\Omega$ . Instead of  $H_0^1(\Omega)$ , we now employ the space

$$V = \left\{ v \in H^1(\Omega) \mid \int_{\Omega} v = 0 \right\}$$

with the operators  $\mathcal{L}$  and  $\mathcal{A}$  defined analogously as above (see (1.3) and (1.4)):

$$\begin{aligned} \langle \mathcal{L}\phi, \psi \rangle &= \int_{\Omega} \nabla\phi \cdot \nabla\psi, \quad \phi, \psi \in V, \\ \langle \mathcal{A}\phi, \psi \rangle &= \int_{\Omega} K\nabla\phi \cdot \nabla\psi, \quad \phi, \psi \in V, \end{aligned}$$

where  $\mathcal{L}$  has a bounded inverse operator; see, e.g., [9, Example 7.2.2, page 117]. For the Neumann problem, the functions  $v_r$ , defined in (2.4), and the solutions  $\phi$  of the wave equation, defined in (2.9), must be modified to

$$v_r - \int_{\Omega} v_r \quad \text{and} \quad \phi - \int_{\Sigma_i} \phi,$$

respectively. The rest will follow in an analogous way to the analysis presented in this paper.

**6. Numerical experiments.** In this section our theoretical results will be illuminated by numerical experiments where the matrices are constructed using FEniCS and the eigenvalues are computed and visualized with MATLAB.<sup>2</sup> If not specified otherwise, we consider the domain  $\Omega \equiv (0, 1) \times (0, 1)$  and a uniform triangular mesh with piecewise linear discretization basis functions is used.

The examples in section 1 concern diagonal positive definite tensors. We first complement this by performing experiments with nondiagonal indefinite tensors. We consider three test problems in the form (1.1) with zero Dirichlet boundary conditions and with the following entries in the symmetric tensor (4.1):

- (P4) :  $k_1(x, y) = 5, k_2(x, y) = -5, k_3(x, y) = 0,$
- (P5) :  $k_1(x, y) = 3, k_2(x, y) = -3, k_3(x, y) = 4,$
- (P6) :  $k_1(x, y) = 3e^{-3(|x-0.5|+|y-0.5|)}, k_2(x, y) = -k_1, k_3(x, y) = 4 \cos\left(\frac{\pi(x+y-1)}{2}\right).$

Using (4.3) gives for problems (P4) and (P5) that  $\kappa_1(x, y) = -5$  and  $\kappa_2(x, y) = 5$ . Furthermore, for problem (P6), formula (4.3) yields

$$\kappa_{1,2}(x, y) = \pm\sqrt{k_1^2 + k_3^2} = \pm\sqrt{9e^{-6(|x-0.5|+|y-0.5|)} + 16 \cos^2\left(\frac{\pi(x+y-1)}{2}\right)},$$

such that  $\kappa_1(\bar{\Omega}) = -\kappa_2(\bar{\Omega}) = [3e^{-3}, 5]$ . As in Figure 1, the spectra visualized in Figure 2 spread over the entire interval (1.8) defined by the nonoverlapping ranges  $\kappa_1(\bar{\Omega})$  and  $\kappa_2(\bar{\Omega})$ . Refining the mesh gives better approximations of the endpoints of the interval  $[-5, 5]$ . The fact that the tensor (4.1) is not diagonal has no qualitative effect on the observed experimental data. We will therefore below only consider diagonal tensors.

<sup>2</sup>FEniCS version 2017.2.0 [1] and MATLAB version 9.5.0 (R2018b).

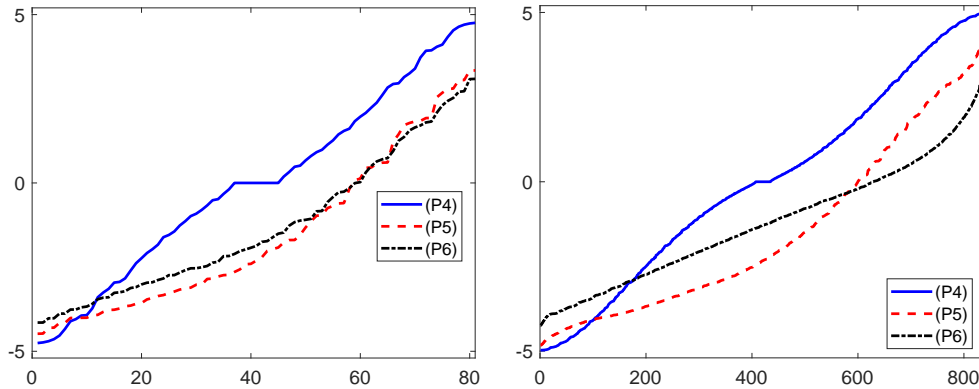


FIG. 2. Spectra of the discretized test problems (P4), (P5), and (P6) for  $N = 81$  (left) and  $N = 841$  (right) degrees of freedom. Horizontal axis: the indices of the increasingly ordered eigenvalues. Vertical axis: the size of the eigenvalues.

The left part of Figure 3 shows numerical results computed with homogeneous Neumann boundary conditions (see section 5). The results with zero Dirichlet boundary conditions are, for comparison, presented in the right part of Figure 3. We consider two test problems with the diagonal tensor (2.1) defined by

$$(6.1) \quad \begin{aligned} \text{(P7)} : \quad & \kappa_1(x, y) = 10 - f(x, y), \quad \kappa_2(x, y) = 4 + f(x, y), \\ \text{(P8)} : \quad & \kappa_1(x, y) = 8 + f(x, y), \quad \kappa_2(x, y) = 6 - f(x, y), \end{aligned}$$

where

$$f(x, y) = 4((x - 0.5)^2 + (y - 0.5)^2)$$

is chosen such that, for both problems,  $\kappa_1(\bar{\Omega}) = [8, 10]$  and  $\kappa_2(\bar{\Omega}) = [4, 6]$ . Note that these intervals do not overlap. The minimum (respectively, maximum) of the interval  $[4, 10]$  is obtained by the function  $\kappa_1(x, y)$  (respectively,  $\kappa_2(x, y)$ ) in the interior of the solution domain for problem (P7), while for problem (P8) the endpoints of this interval are attained on the boundary  $\partial\Omega$ . In the latter case the endpoints of the interval  $[4, 10]$  are better approximated for the problem with Neumann boundary conditions. Similar behavior was also observed for other test cases.

Numerical results for nonconvex domains are presented in Figure 4. We used the diagonal tensor (2.1) with

$$(P9) : \quad \kappa_1(x, y) = 6 - 3e^{-3(|x-0.8|+|y-0.8|)}, \quad \kappa_2(x, y) = 6 + 3e^{-3(|x-0.2|+|y-0.2|)}$$

and the L-shaped domains  $\Omega_1 = (0, 1)^2 \setminus (0, 0.6)^2$  and  $\Omega_2 = (0, 1)^2 \setminus (0.4, 1)^2$ ; see the illustration in the left part of Figure 4. We employed zero Dirichlet boundary conditions. The function  $\kappa_1(x, y)$  (respectively,  $\kappa_2(x, y)$ ) has its minimum (respectively, maximum) at the point  $[0.8, 0.8]$  (respectively,  $[0.2, 0.2]$ ), which is outside the domain  $\Omega_2$  (respectively,  $\Omega_1$ ). As a result, we observe in Figure 4 that the spectra of the discretized problems differ, depending on the ranges of functions  $\kappa_1(x, y)$  and  $\kappa_2(x, y)$  over  $\bar{\Omega}_1$  and  $\bar{\Omega}_2$ .

Finally, we present in Figure 5 numerical results for 3D problems, which are not (yet) supported by rigorous proofs. We consider the unit cube  $\Omega \equiv (0, 1)^3$ , zero Dirichlet boundary conditions, and the diagonal tensor  $K(x, y, z) = \text{diag}(\kappa_1, \kappa_2, \kappa_3)$

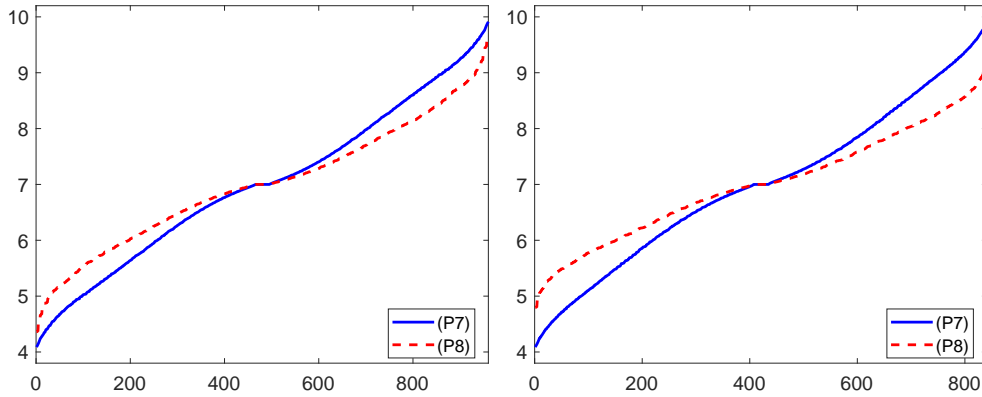


FIG. 3. Spectra of the discretized test problems (P7) and (P8) with zero Neumann boundary conditions (left) and zero Dirichlet boundary conditions (right).

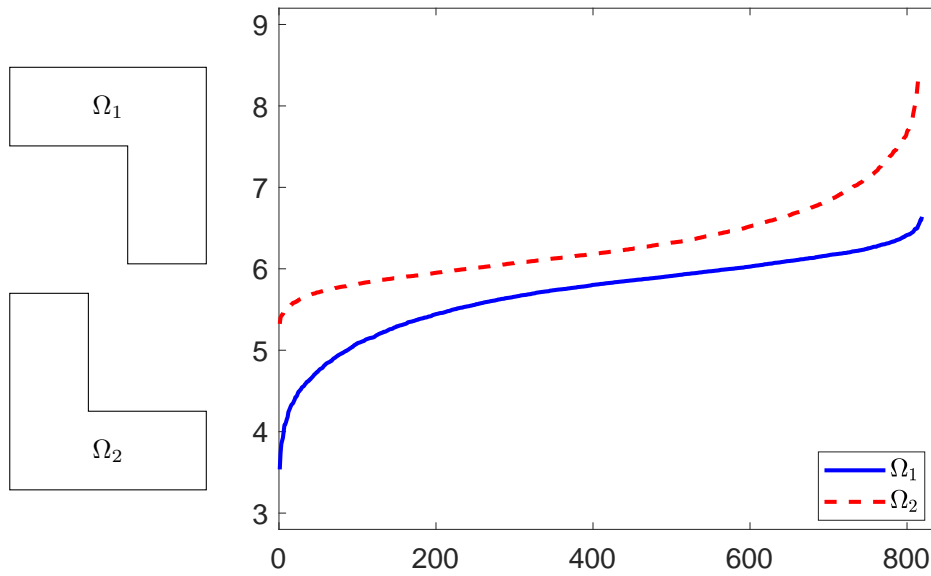


FIG. 4. Left: Illustration of the shapes of the domains  $\Omega_1$  and  $\Omega_2$ . Right: Spectra of the test problem (P9) associated with the domains  $\Omega_1$  and  $\Omega_2$ . The ranges satisfy  $\kappa_1(\overline{\Omega}_1) \subset [3, 6]$  and  $\kappa_2(\overline{\Omega}_1) \subset [6, 7]$  for the domain  $\Omega_1$  and  $\kappa_1(\overline{\Omega}_2) \subset [5, 6]$  and  $\kappa_2(\overline{\Omega}_2) \subset [6, 9]$  for the domain  $\Omega_2$ .

defined as

$$(P10) : \kappa_1 = 1, \quad \kappa_2 = 5.5, \quad \kappa_3 = 10,$$

$$(P11) : \kappa_1 = 1 + \sin^2(x + y + z), \quad \kappa_2 = 5.5 + \cos(\pi xyz), \quad \kappa_3 = 10 - \cos^2(x + y + z),$$

$$(P12) : \kappa_1 = 1 + (x + y + z - 1)^2, \quad \kappa_2 = 4 + xy + z, \quad \kappa_3 = 10 - 2(x + y + z - 1)^2.$$

This choice of test problems follows the same “pattern” as for the introductory experiments presented in section 1: The ranges of the functions  $\kappa_i(x, y, z)$ ,  $i = 1, 2, 3$ , are for (P10) isolated points; they form nonoverlapping intervals for (P11) and overlapping intervals for (P12). As for the 2D test cases, we observe that the spectra of the discretized problems are spread over the entire interval  $[1, 10]$ , irrespective of whether

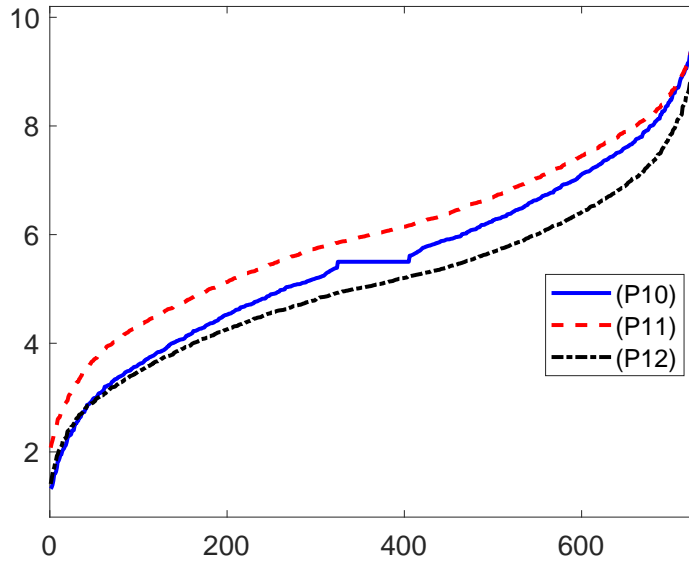


FIG. 5. The spectra of the 3D test problems (P10)–(P12) spread over the entire interval  $[1, 10]$ , while the ranges of the function entries of the diagonal tensors are as follows: Isolated points  $\kappa_1(\bar{\Omega}) = 1$ ,  $\kappa_2(\bar{\Omega}) = 5.5$ , and  $\kappa_3(\bar{\Omega}) = 10$  for (P10); nonoverlapping intervals  $\kappa_1(\bar{\Omega}) = [1, 2]$ ,  $\kappa_2(\bar{\Omega}) = [4.5, 6.5]$ , and  $\kappa_3(\bar{\Omega}) = [9, 10]$  for (P11); and overlapping intervals  $\kappa_1(\bar{\Omega}) = [1, 5]$ ,  $\kappa_2(\bar{\Omega}) = [4, 6]$ , and  $\kappa_3(\bar{\Omega}) = [2, 10]$  for (P12).

the associated ranges overlap or not.

**7. Open problems.** In this paper we have rigorously analyzed 2D problems, and it is an open question whether our main result, Theorem 1.1, also holds in three dimensions or even higher dimensions. Our numerical results indicate that such a generalization is possible, but, e.g., the task of construction functions similar to the  $\{v_r\}$  functions (2.4) will become more involved.

Another important issue is to “translate” our findings to discretized operators. This was accomplished in [4] for uniformly elliptic operators with scalar coefficient functions. That is, [4] contains discrete versions of the results published in [10] and further develops towards approximating locally the individual eigenvalues. The techniques employed in [4] can be generalized to analyze discretized second order differential operators with indefinite tensors. Such a development is, however, out of the scope of this paper. An interesting question concerns the distribution of the eigenvalues: For discretized operators, are the eigenvalues evenly distributed in the interval (1.8)? Our numerical experiments suggest that the answer may be positive. We will return to this question elsewhere.

**Appendix A. Technical details about the inequalities (2.5) in the proof of Lemma 2.1.** We want to prove that, for sufficiently small  $r > 0$ ,

$$(A.1) \quad 4 - 4r \leq \|(v_r)_x\|_{L^2(\Omega)}^2 \leq 4,$$

$$(A.2) \quad \|(v_r)_y\|_{L^2(\Omega)}^2 \leq 4r,$$

where  $v_r(x, y)$  is defined on  $R_r$  by (2.3) and (2.4). Without loss of generality, we

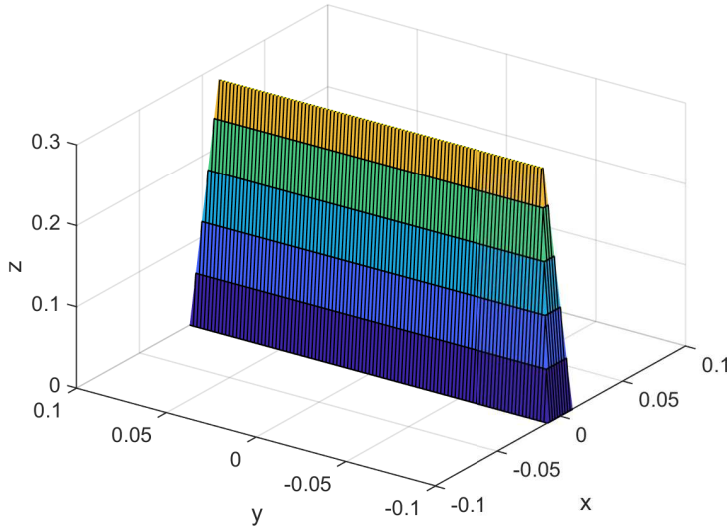


FIG. 6. The function  $v_r$  centered at the point  $(0,0)$  with  $r = 0.1$ .

consider the case  $(x_0, y_0) = (0, 0)$ . Then  $R_r = [-r^2, r^2] \times [-r, r]$  and

$$v_r(x, y) = \sqrt{r} \min \left\{ 1 - \frac{|x|}{r^2}, \frac{1}{r} - \frac{|y|}{r^2} \right\} \quad \text{for } (x, y) \in R_r,$$

with  $v_r(x, y) = 0$  elsewhere; see Figure 6.

For any  $0 < r < 1$ , the partial derivatives of  $v_r(x, y)$  are not defined at the boundary  $\partial R_r$  of  $R_r$ , at the set of points  $\{(x, y) \in R_r : |y| - |x| = r - r^2\}$ , and at the set of points  $\{(x, y) : x = 0, |y| < r - r^2\}$  where  $v_r(x, y)$  reaches its maximum; see the edges of  $\{v_r(R_r)\}$  in Figure 6. Simple computations yield that within  $R_r$

$$\begin{aligned} |\partial_x v_r(x, y)|^2 &= 0, & |\partial_y v_r(x, y)|^2 &= \frac{1}{r^3} & \text{for } |y| - |x| > r - r^2, & (x, y) \notin \partial R_r, \\ |\partial_x v_r(x, y)|^2 &= \frac{1}{r^3}, & |\partial_y v_r(x, y)|^2 &= 0 & \text{for } x \neq 0, |y| - |x| < r - r^2, & (x, y) \notin \partial R_r. \end{aligned}$$

The upper bound in (A.1) thus holds because

$$(A.3) \quad \|(v_r)_x\|_{L^2(\Omega)}^2 = \int_{R_r} |\partial_x v_r(x, y)|^2 \leq \int_{R_r} \frac{1}{r^3} = \frac{2r^2 \cdot 2r}{r^3} = 4.$$

Moreover, denoting

$$P_r = \{(x, y) : x \neq 0, |x| < r^2, |y| < r - r^2\},$$

we have

$$|\partial_x v_r(x, y)|^2 = \frac{1}{r^3}, \quad |\partial_y v_r(x, y)|^2 = 0 \quad \text{for } (x, y) \in P_r.$$

Thus  $\|(v_r)_x\|_{L^2(\Omega)}^2$  and  $\|(v_r)_y\|_{L^2(\Omega)}^2$  can be bounded as follows:

$$\begin{aligned} \|(v_r)_x\|_{L^2(\Omega)}^2 &= \int_{R_r} |\partial_x v_r(x, y)|^2 \geq \int_{P_r} |\partial_x v_r(x, y)|^2 = \int_{P_r} \frac{1}{r^3} = \frac{2r^2 \cdot 2(r-r^2)}{r^3} = 4 - 4r, \\ \|(v_r)_y\|_{L^2(\Omega)}^2 &= \int_{R_r} |\partial_y v_r(x, y)|^2 = \int_{R_r \setminus P_r} |\partial_y v_r(x, y)|^2 \leq \int_{R_r \setminus P_r} \frac{1}{r^3} = \frac{2r^2 \cdot 2r^2}{r^3} = 4r, \end{aligned}$$

which completes the proof.

**Appendix B. Technical details about the bound (2.18) in the proof of Lemma 2.3.** Assume that  $\mathcal{L}^{-1}\mathcal{A} - \lambda\mathcal{I}$  has a bounded inverse. We will show that, for sufficiently small  $l > 0$ ,

$$(B.1) \quad \|(\mathcal{L}^{-1}\mathcal{A} - \lambda\mathcal{I})^{-1}\mathcal{L}^{-1}(\mathcal{A}_l - \mathcal{A})\|_{\mathcal{L}} \leq \|(\mathcal{L}^{-1}\mathcal{A} - \lambda\mathcal{I})^{-1}\|_{\mathcal{L}}\|\mathcal{L}^{-1}(\mathcal{A}_l - \mathcal{A})\|_{\mathcal{L}} < 1.$$

The operator norm

$$\|\mathcal{L}^{-1}(\mathcal{A}_l - \mathcal{A})\|_{\mathcal{L}} \equiv \sup_{u \in H_0^1(\Omega)} \frac{\|\mathcal{L}^{-1}(\mathcal{A}_l - \mathcal{A})u\|_{\mathcal{L}}}{\|u\|_{\mathcal{L}}}$$

can be expressed as (see, e.g. [2, Theorem 4.1–3])

$$(B.2) \quad \|\mathcal{L}^{-1}(\mathcal{A}_l - \mathcal{A})\|_{\mathcal{L}} = \sup_{u, v \in H_0^1(\Omega)} \frac{|(\mathcal{L}^{-1}(\mathcal{A}_l - \mathcal{A})u, v)_{\mathcal{L}}|}{\|u\|_{\mathcal{L}}\|v\|_{\mathcal{L}}}.$$

Using

$$\begin{aligned} |(\mathcal{L}^{-1}(\mathcal{A}_l - \mathcal{A})u, v)_{\mathcal{L}}| &= |((\mathcal{A}_l - \mathcal{A})u, v)| \\ &= \left| \int_{S_l} (K(x_0, y_0) - K(x, y)) \nabla u \cdot \nabla v \right| \\ &\leq \int_{S_l} \|K(x_0, y_0) - K(x, y)\| |\nabla u| \cdot |\nabla v| \\ &\leq \sup_{(x, y) \in S_l} \|K(x_0, y_0) - K(x, y)\| \|u\|_{\mathcal{L}} \|v\|_{\mathcal{L}}, \end{aligned}$$

we get the bound

$$\|\mathcal{L}^{-1}(\mathcal{A}_l - \mathcal{A})\|_{\mathcal{L}} \leq \sup_{(x, y) \in S_l} \|K(x_0, y_0) - K(x, y)\|.$$

Since  $\|(\mathcal{L}^{-1}\mathcal{A} - \lambda\mathcal{I})^{-1}\|_{\mathcal{L}}$  is bounded, the continuity of  $K(x, y)$  at the point  $(x_0, y_0)$  ensures that  $l$  can be chosen such that (B.1) holds.

#### REFERENCES

- [1] M. S. ALNÆS, J. BLECHTA, J. HAKE, A. JOHANSSON, B. KEHLET, A. LOGG, C. RICHARDSON, J. RING, M. E. ROGNES, AND G. N. WELLS, *The FEniCS project version 1.5*, Arch. Numer. Software, 3 (2015), pp. 9–23, <https://doi.org/10.11588/ans.2015.100.20553>.
- [2] P. G. CIARLET, *Linear and Nonlinear Functional Analysis with Applications*, SIAM, Philadelphia, 2013.
- [3] A. F. FRIEDMAN, *Foundations of Modern Analysis*, Dover, New York, 1982.
- [4] T. GERGELITS, K.-A. MARDAL, B. F. NIELSEN, AND Z. STRAKOŠ, *Laplacian preconditioning of elliptic PDEs: Localization of the eigenvalues of the discretized operator*, SIAM J. Numer. Anal., 57 (2019), pp. 1369–1394, <https://doi.org/10.1137/18M1212458>.
- [5] A. GREENBAUM, *Behavior of slightly perturbed Lanczos and conjugate-gradient recurrences*, Linear Algebra Appl., 113 (1989), pp. 7–63, [https://doi.org/10.1016/0024-3795\(89\)90285-1](https://doi.org/10.1016/0024-3795(89)90285-1).
- [6] C. HOFREITHER AND S. TAKACS, *Robust multigrid for isogeometric analysis based on stable splittings of spline spaces*, SIAM J. Numer. Anal., 55 (2017), pp. 2004–2024, <https://doi.org/10.1137/16M1085425>.
- [7] C. HOFREITHER, S. TAKACS, AND W. ZULEHNER, *A robust multigrid method for isogeometric analysis in two dimensions using boundary correction*, Comput. Methods Appl. Mech. Engrg., 316 (2017), pp. 22–42.

- [8] J. LIESEN AND Z. STRAKOŠ, *Krylov Subspace Methods: Principles and Analysis*, Numer. Math. Sci. Comput., Oxford University Press, Oxford, UK, 2013.
- [9] J. T. MARTI, *Introduction to Sobolev Spaces and Finite Element Solution of Elliptic Boundary Value Problems*, Academic Press, London, 1986.
- [10] B. F. NIELSEN, A. TVEITO, AND W. HACKBUSCH, *Preconditioning by inverting the Laplacian; an analysis of the eigenvalues*, IMA J. Numer. Anal., 29 (2009), pp. 24–42.
- [11] M. RENARDY AND R. C. ROGERS, *An Introduction to Partial Differential Equations*, Springer-Verlag, New York, 1993.
- [12] G. SANGALLI AND M. TANI, *Isogeometric preconditioners based on fast solvers for the Sylvester equation*, SIAM J. Sci. Comput., 38 (2016), pp. A3644–A3671, <https://doi.org/10.1137/16M1062788>.