

# Spectral analysis and spectral symbol of matrices in isogeometric Galerkin methods

Carlo Garoni<sup>a,b</sup>, Carla Manni<sup>a</sup>, Stefano Serra-Capizzano<sup>b,c</sup>, Debora Sesana<sup>b</sup>,  
Hendrik Speleers<sup>a</sup>

<sup>a</sup>University of Roma 'Tor Vergata', Department of Mathematics, Via della Ricerca Scientifica, 00133 Roma, Italy.  
Email: garoni@mat.uniroma2.it, manni@mat.uniroma2.it, speleers@mat.uniroma2.it.

<sup>b</sup>University of Insubria, Department of Science and High Technology, Via Valleggio 11, 22100 Como, Italy.  
Email: carlo.garoni@uninsubria.it, stefano.serrac@uninsubria.it, debora.sesana@uninsubria.it.

<sup>c</sup>Division of Scientific Computing, Department of Information Technology, Uppsala University,  
Box 337, SE-751 05 Uppsala, Sweden. Email: stefano.serra@it.uu.se.

January 29, 2015

## Abstract

A linear full elliptic second order Partial Differential Equation (PDE), defined on a  $d$ -dimensional domain  $\Omega$ , is approximated by the isogeometric Galerkin method based on uniform tensor-product B-splines of degrees  $(p_1, \dots, p_d)$ . The considered approximation process leads to a  $d$ -level stiffness matrix, banded in a multilevel sense. This matrix is close to a  $d$ -level Toeplitz structure when the PDE coefficients are constant and the physical domain  $\Omega$  is just the hypercube  $(0, 1)^d$  without using any geometry map. In such a simplified case, a detailed spectral analysis of the stiffness matrices has been carried out in a previous work. In this paper, we complete the picture by considering non-constant PDE coefficients and an arbitrary domain  $\Omega$ , parameterized with a non-trivial geometry map. We compute and study the spectral symbol of the related stiffness matrices. This symbol describes the asymptotic eigenvalue distribution when the fineness parameters tend to zero (so that the matrix-size tends to infinity). The mathematical technique used for computing the symbol is based on the theory of Generalized Locally Toeplitz (GLT) sequences.

*Keywords:* spectral distribution, symbol, Galerkin method, B-splines, Isogeometric Analysis.

*2010 MSC:* 15A18, 15B05, 41A15, 15A69, 65N30.

## 1 Introduction

Isogeometric Analysis (IgA) is a paradigm for the analysis of problems governed by Partial Differential Equations (PDEs); see [5]. Its goal is to improve the connection between numerical simulation and Computer Aided Design (CAD) systems. In its original formulation, the main idea in IgA is to use directly the geometry provided by CAD systems and to approximate the unknown solutions of differential equations by the same type of functions. Tensor-product B-splines and their rational extension, the so-called NURBS, are the dominant technology in CAD systems used in engineering, and thus also in IgA.

---

This work was partially supported by INdAM-GNCS Gruppo Nazionale per il Calcolo Scientifico, by the MIUR 'Futuro in Ricerca 2013' Programme through the project DREAMS, and by the Program 'Becoming the Number One – Sweden (2014)' of the Knut and Alice Wallenberg Foundation.

In this paper, we consider the following linear full elliptic second order PDE, with non-constant coefficients and homogeneous Dirichlet boundary conditions:

$$\begin{cases} -\nabla \cdot K \nabla u + \alpha \cdot \nabla u + \gamma u = f, & \text{in } \Omega, \\ u = 0, & \text{on } \partial\Omega, \end{cases} \quad (1.1)$$

where  $\Omega$  is a bounded open domain in  $\mathbb{R}^d$  with Lipschitz boundary,  $K : \Omega \rightarrow \mathbb{R}^{d \times d}$  is a Symmetric Positive Definite (SPD) matrix of functions in  $L^\infty(\Omega)$ ,  $\alpha : \Omega \rightarrow \mathbb{R}^d$  is a vector of functions in  $L^\infty(\Omega)$ ,  $\gamma \in L^\infty(\Omega)$ ,  $\gamma \geq 0$  and  $f \in L^2(\Omega)$ . We focus on the isogeometric Galerkin approximation of (1.1) using a general geometry map  $\mathbf{G} : [0, 1]^d \rightarrow \bar{\Omega}$  and uniform tensor-product B-splines of arbitrary degrees  $\mathbf{p} := (p_1, \dots, p_d)$ .

This paper is devoted to the study of the asymptotic spectral (and singular value) distribution of the resulting Galerkin B-spline IgA stiffness matrices, when the fineness parameters tend to zero so that the related matrix-size tends to infinity. Our results extend those obtained in [11] and generalized in [10, Chapter 4], which address the simplified case where  $K$  is the identity matrix,  $\Omega = (0, 1)^d$ , and  $\mathbf{G}$  is the identity map. More in detail, in this paper we prove the following:

- a) a spectral distribution exists and is compactly described by a symbol  $f$ ;
- b) the symbol  $f$  has a canonical structure incorporating:
  - b1) the approximation technique, identified by a finite set of polynomials in the Fourier variables  $\boldsymbol{\theta} := (\theta_1, \dots, \theta_d) \in [-\pi, \pi]^d$ ;
  - b2) the geometry, identified by the map  $\mathbf{G}$  in the parametric variables  $\hat{\mathbf{x}} := (\hat{x}_1, \dots, \hat{x}_d)$  defined on the reference domain  $[0, 1]^d$ ;
  - b3) the coefficients of the higher-order operator of the PDE, namely  $K$ , in the physical variables  $\mathbf{x} := (x_1, \dots, x_d)$  defined on the physical domain  $\Omega$ ;
- c) the symbol  $f$  is the same as in the isogeometric collocation setting [8], up to a determinant factor  $|\det(J_{\mathbf{G}})|$ , being  $J_{\mathbf{G}}$  the Jacobian matrix of  $\mathbf{G}$ ;
- d) when  $K$  is the identity matrix,  $\Omega = (0, 1)^d$  and  $\mathbf{G}$  is the identity map, the symbol  $f$  reduces, as expected, to the one in [10, Chapter 4].

The picture described in item b) is intrinsic to the approximation of PDEs by any local method, such as Finite Differences (FDs) and Finite Elements (FEs). Actually, the formal structure of the symbol is essentially the same when considering different techniques to approximate the same problem; see [1, 20, 21] and references therein, with particular attention to [20, Section 2] and [21, Question 3.1]. Moreover, the appearance of the determinant factor mentioned in item c) is expected. Indeed, this behavior was already observed when passing from FDs (whose philosophy is analogous in the collocation framework) to FEs in the Galerkin context. Note that the determinant factor helps in keeping the IgA Galerkin matrices less ill-conditioned than the corresponding IgA collocation matrices, even when the map  $\mathbf{G}$  is (nearly) singular.

Although the formal structure of the symbol is shared by different approximation techniques, some of its analytic features are not so common. For instance, if one of the  $\mathbf{p}$  parameters, say  $p_i$ , becomes large, then the symbol  $f$  shows ‘numerical zeros’ at the points  $(\hat{\mathbf{x}}, \boldsymbol{\theta}) \in [0, 1]^d \times [-\pi, \pi]^d$  where  $\theta_i = \pi$ . More precisely, if  $\theta_i = \pi$ , the value  $f(\hat{\mathbf{x}}, \boldsymbol{\theta})$  converges to 0 exponentially when  $p_i \rightarrow \infty$ . The latter information implies that small eigenvalues appear related to high frequency eigenvectors, and this non-canonical source of ill-conditioning is responsible for the slowdown of all the standard multigrid and preconditioning techniques when one of the  $p_i$  grows. On the other hand, quite recently, we have found a way of exploiting the spectral information provided by the symbol  $f$  for designing algorithms with convergence speed independent of the fineness parameters and also of the approximation parameters  $\mathbf{p}$ ; see [6, 7, 9].

We would like to emphasize that, besides the identification of the symbol  $f$  for the Galerkin B-spline IgA stiffness matrices, the other important aspect of this paper is the mathematical technique used in our derivation. As explained in Section 4, this technique is quite general and can be also applied in other contexts than Galerkin B-spline IgA. It consists of the following mathematical tools. We use the theory of separable (multilevel) Locally Toeplitz (sLT) sequences and the theory of Generalized Locally Toeplitz (GLT) sequences, which go back to the pioneering work by Tilli [23] and are developed in [20, 21]. Implicitly, we also use the concept of approximating class of sequences (a.c.s.), which was introduced in [18] and allows one to derive the singular value and eigenvalue distribution of a complicated sequence of matrices (matrix-sequence) starting from those of simpler matrix-sequences; see [14, 18]. Finally, we exploit general results, contained in [15] and generalized in [12], which allow one to determine the spectral distribution of arbitrary (non-Hermitian) perturbed versions of sequences formed by Hermitian matrices, under certain conditions on the perturbation matrices.

Another way to obtain the main results of this paper could have been a comparison technique between the Galerkin B-spline IgA stiffness matrices considered herein and the B-spline IgA collocation matrices analyzed in [8] and [10, Chapter 5]. However, we preferred the approach discussed in the previous paragraph, due to its intrinsic generality.

The paper is organized as follows. In Section 2 we introduce the notation and definitions used throughout the paper; we also report some basic results. In Section 3 we describe the isogeometric Galerkin approximation based on uniform tensor-product B-splines of degrees  $\mathbf{p}$  of the full elliptic problem (1.1). Sections 4 and 5 contain our main results: the computation of the spectral and singular value distribution of the Galerkin B-spline IgA stiffness matrices, the identification of the corresponding symbol  $f$ , and the study of its properties. We end in Section 6 with some concluding remarks.

## 2 Preliminaries

### 2.1 Multi-index notation

Throughout this paper, we will use the multi-index notation, expounded by Tyrtysnikov in [24, Section 6]. When discretizing a linear PDE over a  $d$ -dimensional domain by means of some numerical method, the resulting discretization matrices show a  $d$ -level structure; see [24, Section 6] for the corresponding definition. The multi-index notation is a powerful tool that allows one to give a compact expression of these matrices, treating the dimensionality parameter  $d$  as any other parameter involved in the discretization process. In this way, the dependency of the matrix structure on  $d$  is highlighted and a compact presentation is made possible.

A multi-index  $\mathbf{m} \in \mathbb{Z}^d$ , also called a  $d$ -index, is simply a (row) vector in  $\mathbb{Z}^d$  and its components are denoted by  $m_1, \dots, m_d$ . We indicate by  $\mathbf{0}$ ,  $\mathbf{1}$ ,  $\mathbf{2}$  the vectors consisting of all zeros, all ones, all twos, respectively (their size will be clear from the context). For any  $d$ -index  $\mathbf{m}$ , we set  $N(\mathbf{m}) := \prod_{i=1}^d m_i$  and we write  $\mathbf{m} \rightarrow \infty$  to indicate that  $\min_{i=1, \dots, d} m_i \rightarrow \infty$ . Inequalities between multi-indices must be interpreted in the componentwise sense. For example,  $\mathbf{j} \leq \mathbf{k}$  means that  $j_i \leq k_i$  for every  $i$ . If  $\mathbf{j}, \mathbf{k}$  are  $d$ -indices such that  $\mathbf{j} \leq \mathbf{k}$ , the multi-index range  $\mathbf{j}, \dots, \mathbf{k}$  is the set  $\{\mathbf{i} \in \mathbb{Z}^d : \mathbf{j} \leq \mathbf{i} \leq \mathbf{k}\}$ . We assume for this set the standard lexicographic ordering:

$$\left[ \dots \left[ \left[ (i_1, \dots, i_d) \right]_{i_d=j_d, \dots, k_d} \right]_{i_{d-1}=j_{d-1}, \dots, k_{d-1}} \dots \right]_{i_1=j_1, \dots, k_1}. \quad (2.1)$$

For instance, in the case  $d = 2$ , this ordering is

$$(j_1, j_2), (j_1, j_2 + 1), \dots, (j_1, k_2), (j_1 + 1, j_2), (j_1 + 1, j_2 + 1), \dots, (j_1 + 1, k_2), \\ \dots \dots \dots, (k_1, j_2), (k_1, j_2 + 1), \dots, (k_1, k_2).$$

When a  $d$ -index  $i$  varies in a multi-index range  $\mathbf{j}, \dots, \mathbf{k}$  (this is sometimes written as  $i = \mathbf{j}, \dots, \mathbf{k}$ ), it is always assumed that  $i$  varies from  $\mathbf{j}$  to  $\mathbf{k}$  following the specific ordering (2.1). In particular, if  $\mathbf{m} \in \mathbb{N}^d$  and  $\mathbf{x} = [x_i]_{i=1}^{\mathbf{m}}$ , then  $\mathbf{x}$  is a vector of length  $N(\mathbf{m})$  whose components  $x_i$ ,  $i = 1, \dots, \mathbf{m}$ , are ordered in accordance with (2.1): the first component is  $x_1 = x_{(1, \dots, 1, 1)}$ , the second component is  $x_{(1, \dots, 1, 2)}$ , and so on until the last component, which is  $x_{\mathbf{m}} = x_{(m_1, \dots, m_d)}$ . Similarly, if  $X = [x_{ij}]_{i,j=1}^{\mathbf{m}}$ , then  $X$  is a  $N(\mathbf{m}) \times N(\mathbf{m})$  matrix whose entries are indexed by two  $d$ -indices  $\mathbf{i}, \mathbf{j}$ , both varying over the multi-index range  $1, \dots, \mathbf{m}$  according to (2.1). Operations involving multi-indices that do not have a meaning when considering multi-indices as normal vectors must always be interpreted in the componentwise sense. For example,  $\mathbf{j}\mathbf{k} := (j_1 k_1, \dots, j_d k_d)$ ,  $\mathbf{j}/\mathbf{k} := (j_1/k_1, \dots, j_d/k_d)$ , etc.

## 2.2 Preliminaries on matrix analysis

For all  $X \in \mathbb{C}^{m \times m}$ , the singular values of  $X$  are denoted by  $\sigma_j(X)$ ,  $j = 1, \dots, m$ , and the eigenvalues of  $X$  by  $\lambda_j(X)$ ,  $j = 1, \dots, m$ . If  $X, Y \in \mathbb{C}^{m \times m}$ , the notation  $X \geq Y$  (resp.  $X > Y$ ) means that  $X, Y$  are Hermitian and  $X - Y$  is positive semi-definite (resp. positive definite). The  $\infty$ -norm and the spectral norm (2-norm) of both vectors and matrices are denoted by  $\|\cdot\|_\infty$  and  $\|\cdot\|$ , respectively. We recall that, for all  $X \in \mathbb{C}^{m \times m}$ ,

$$\|X\| \leq \sqrt{\|X\|_\infty \|X^T\|_\infty}, \quad \forall X \in \mathbb{C}^{m \times m}. \quad (2.2)$$

For  $X \in \mathbb{C}^{m \times m}$ , we denote by  $\|X\|_1$  the trace norm of  $X$ , i.e., the sum of all the singular values of  $X$ . The trace norm, also called Schatten 1-norm, and the other Schatten  $p$ -norms are studied in [2]. Since  $\text{rank}(X)$  is the number of non-zero singular values of  $X$  and  $\|X\|$  equals the largest singular value of  $X$ , we have

$$\|X\|_1 \leq \text{rank}(X) \|X\| \leq m \|X\|, \quad \forall X \in \mathbb{C}^{m \times m}. \quad (2.3)$$

If  $X, Y$  are matrices of any dimension, say  $X \in \mathbb{C}^{\ell \times m}$  and  $Y \in \mathbb{C}^{q \times r}$ , the tensor (Kronecker) product  $X \otimes Y$  is defined as

$$X \otimes Y := [x_{ij} Y]_{\substack{i=1, \dots, \ell \\ j=1, \dots, m}} = \begin{bmatrix} x_{11} Y & \cdots & x_{1m} Y \\ \vdots & & \vdots \\ x_{\ell 1} Y & \cdots & x_{\ell m} Y \end{bmatrix} \in \mathbb{C}^{\ell q \times m r}.$$

Tensor products possess a lot of nice algebraic properties. One of them is the associativity, which allows one to omit parentheses in expressions like  $X_1 \otimes X_2 \otimes \cdots \otimes X_d$ . Another property is the following [10, Section 1.2.1]: let  $X_1, \dots, X_d, Y_1, \dots, Y_d$  be matrices such that  $X_i, Y_i \in \mathbb{C}^{m_i \times m_i}$  for all  $i = 1, \dots, d$ , then

$$\text{rank}(X_1 \otimes \cdots \otimes X_d - Y_1 \otimes \cdots \otimes Y_d) \leq N(\mathbf{m}) \sum_{i=1}^d \frac{\text{rank}(X_i - Y_i)}{m_i}, \quad (2.4)$$

where  $\mathbf{m} := (m_1, \dots, m_d)$ .

## 2.3 Spectral distribution and spectral symbol

We denote by  $\mu_d$  the Lebesgue measure in  $\mathbb{R}^d$  and by  $C_c(\mathbb{C})$  the space of continuous functions  $F : \mathbb{C} \rightarrow \mathbb{C}$  with bounded support.

**Definition 2.1.** Let  $\{X_n\}_n$  be a sequence of matrices, with  $X_n$  of size  $d_n$  tending to infinity, and let  $f : D \rightarrow \mathbb{C}$  be a measurable function defined on the measurable set  $D \subset \mathbb{R}^d$ , with  $0 < \mu_d(D) < \infty$ . We say that  $\{X_n\}_n$  is distributed like  $f$  in the sense of the singular values and we write  $\{X_n\}_n \sim_\sigma f$ , if

$$\lim_{n \rightarrow \infty} \frac{1}{d_n} \sum_{j=1}^{d_n} F(\sigma_j(X_n)) = \frac{1}{\mu_d(D)} \int_D F(|f(x_1, \dots, x_d)|) dx_1 \cdots dx_d, \quad \forall F \in C_c(\mathbb{C}).$$

In this case,  $f$  is referred to as the singular value symbol of  $\{X_n\}_n$ .

Similarly, we say that  $\{X_n\}_n$  is distributed like  $f$  in the sense of the eigenvalues and we write  $\{X_n\}_n \sim_\lambda f$ , if

$$\lim_{n \rightarrow \infty} \frac{1}{d_n} \sum_{j=1}^{d_n} F(\lambda_j(X_n)) = \frac{1}{\mu_d(D)} \int_D F(f(x_1, \dots, x_d)) dx_1 \cdots dx_d, \quad \forall F \in C_c(\mathbb{C}).$$

In this case,  $f$  is referred to as the eigenvalue (or spectral) symbol of  $\{X_n\}_n$ .

We refer to [6, Remark 3.1] for the informal meaning behind Definition 2.1.

## 2.4 Toeplitz and diagonal sampling matrices

Given  $\mathbf{m} \in \mathbb{N}^d$  and a function  $f : [-\pi, \pi]^d \rightarrow \mathbb{C}$  belonging to  $L^1([-\pi, \pi]^d)$ , the  $d$ -level Toeplitz matrix  $T_{\mathbf{m}}(f)$  associated with  $f$  is defined as follows [4, 24]:

$$T_{\mathbf{m}}(f) := [\hat{f}_{i-j}]_{i,j=1}^{\mathbf{m}},$$

where

$$\hat{f}_{\mathbf{k}} := \frac{1}{(2\pi)^d} \int_{[-\pi, \pi]^d} f(\boldsymbol{\theta}) e^{-i\mathbf{k} \cdot \boldsymbol{\theta}} d\boldsymbol{\theta}, \quad \mathbf{k} \in \mathbb{Z}^d,$$

are the Fourier coefficients of  $f$  and  $\mathbf{k} \cdot \boldsymbol{\theta} := \sum_{i=1}^d k_i \theta_i$ . The function  $f$  is referred to as the generating function of the Toeplitz family  $\{T_{\mathbf{m}}(f)\}_{\mathbf{m} \in \mathbb{N}^d}$ . The Toeplitz operator  $T_{\mathbf{m}}(\cdot) : L^1([-\pi, \pi]^d) \rightarrow \mathbb{C}^{N(\mathbf{m}) \times N(\mathbf{m})}$  is linear, so for every  $a, b \in \mathbb{C}$  and every  $f, g \in L^1([-\pi, \pi]^d)$  we have

$$T_{\mathbf{m}}(af + bg) = aT_{\mathbf{m}}(f) + bT_{\mathbf{m}}(g).$$

If  $f_i : E_i \subseteq \mathbb{R}^{\ell_i} \rightarrow \mathbb{C}$ ,  $i = 1, \dots, d$ , the tensor-product function  $f_1 \otimes \cdots \otimes f_d : E_1 \times \cdots \times E_d \rightarrow \mathbb{C}$  is defined as

$$(f_1 \otimes \cdots \otimes f_d)(\mathbf{x}_1, \dots, \mathbf{x}_d) := f_1(\mathbf{x}_1) \cdots f_d(\mathbf{x}_d).$$

In the case where  $f_i \in L^1(E_i)$  for all  $i$ , we have  $f_1 \otimes \cdots \otimes f_d \in L^1(E_1 \times \cdots \times E_d)$ . The next result relates tensor products and Toeplitz matrices; see, e.g., [10, Lemma 1.8].

**Lemma 2.2.** *Given  $f_1, \dots, f_d \in L^1([-\pi, \pi])$  and  $\mathbf{m} = (m_1, \dots, m_d) \in \mathbb{N}^d$ , we have*

$$T_{\mathbf{m}}(f_1 \otimes \cdots \otimes f_d) = T_{m_1}(f_1) \otimes \cdots \otimes T_{m_d}(f_d).$$

Given  $\mathbf{m} \in \mathbb{N}^d$  and  $a : [0, 1]^d \rightarrow \mathbb{C}$ , the  $d$ -level diagonal sampling matrix  $D_{\mathbf{m}}(a)$  associated with  $a$  is the  $N(\mathbf{m}) \times N(\mathbf{m})$  matrix defined by

$$D_{\mathbf{m}}(a) := \text{diag}_{j=1, \dots, m} a\left(\frac{\mathbf{j}}{\mathbf{m}}\right), \quad (2.5)$$

where  $\mathbf{j}$  varies from 1 to  $\mathbf{m}$  following the lexicographic ordering (2.1), as explained in Section 2.1.

## 3 Isogeometric Galerkin B-spline approximation

### 3.1 Isogeometric Galerkin methods

The weak form of (1.1) consists in finding  $u \in H_0^1(\Omega)$  such that

$$a(u, v) = F(v), \quad \forall v \in H_0^1(\Omega),$$

where  $a(u, v) := \int_{\Omega} ((\nabla u)^T K \nabla v + (\nabla u)^T \alpha v + \gamma uv)$  and  $F(v) := \int_{\Omega} f v$ . In the standard Galerkin method, we look for an approximation  $u_{\mathcal{W}}$  of  $u$  by choosing a finite dimensional approximation space  $\mathcal{W} \subset H_0^1(\Omega)$  and by solving the following (Galerkin) problem: find  $u_{\mathcal{W}} \in \mathcal{W}$  such that

$$a(u_{\mathcal{W}}, v) = F(v), \quad \forall v \in \mathcal{W}.$$

If  $\{\varphi_1, \dots, \varphi_N\}$  is a basis of  $\mathcal{W}$ , then we can write  $u_{\mathcal{W}} = \sum_{j=1}^N u_j \varphi_j$  for a unique vector  $\mathbf{u} := (u_1, \dots, u_N)^T$ , and the computation of  $u_{\mathcal{W}}$  is equivalent to solving the linear system

$$\mathbf{A} \mathbf{u} = \mathbf{f},$$

where

$$A := [a(\varphi_j, \varphi_i)]_{i,j=1}^N = \left[ \int_{\Omega} ((\nabla \varphi_j)^T K \nabla \varphi_i + (\nabla \varphi_j)^T \alpha \varphi_i + \gamma \varphi_j \varphi_i) \right]_{i,j=1}^N \quad (3.1)$$

is the stiffness matrix and  $\mathbf{f} := [F(\varphi_i)]_{i=1}^N$ .

Now, suppose that the physical domain  $\Omega$  can be described by a global geometry function  $\mathbf{G} : [0, 1]^d \rightarrow \overline{\Omega}$ , which is invertible and satisfies  $\mathbf{G}(\partial([0, 1]^d)) = \partial \overline{\Omega}$ . Let

$$\{\hat{\varphi}_1, \dots, \hat{\varphi}_N\} \quad (3.2)$$

be a set of basis functions defined on the reference (parametric) domain  $[0, 1]^d$  and vanishing on the boundary  $\partial([0, 1]^d)$ . In the Galerkin IgA approach, we find an approximation  $u_{\mathcal{W}}$  of  $u$  by using the standard Galerkin method, in which the approximation space is chosen as  $\mathcal{W} := \langle \varphi_i : i = 1, \dots, N \rangle$ , with

$$\varphi_i(\mathbf{x}) := \hat{\varphi}_i(\mathbf{G}^{-1}(\mathbf{x})) = \hat{\varphi}_i(\hat{\mathbf{x}}), \quad \mathbf{x} = \mathbf{G}(\hat{\mathbf{x}}). \quad (3.3)$$

The resulting stiffness matrix  $A$  is given by (3.1), with the basis functions  $\varphi_i$  defined in (3.3). Assuming that  $\mathbf{G}$  and  $\hat{\varphi}_i$ ,  $i = 1, \dots, N$ , are sufficiently regular, we can apply standard differential calculus to obtain the following expression for  $A$  in terms of  $\mathbf{G}$  and  $\hat{\varphi}_i$ ,  $i = 1, \dots, N$ :

$$A = \left[ \int_{[0,1]^d} ((\nabla \hat{\varphi}_j)^T K_{\mathbf{G}} \nabla \hat{\varphi}_i + (\nabla \hat{\varphi}_j)^T (J_{\mathbf{G}})^{-1} \alpha(\mathbf{G}) \hat{\varphi}_i + \gamma(\mathbf{G}) \hat{\varphi}_j \hat{\varphi}_i) |\det(J_{\mathbf{G}})| \right]_{i,j=1}^N, \quad (3.4)$$

where

$$K_{\mathbf{G}} := (J_{\mathbf{G}})^{-1} K(\mathbf{G}) (J_{\mathbf{G}})^{-T} \quad (3.5)$$

and  $J_{\mathbf{G}}$  is the Jacobian matrix of  $\mathbf{G}$ , i.e.,

$$J_{\mathbf{G}} := \left[ \frac{\partial G_i}{\partial \hat{x}_j} \right]_{i,j=1}^d = \left[ \frac{\partial x_i}{\partial \hat{x}_j} \right]_{i,j=1}^d.$$

In the context of IgA, the geometry map  $\mathbf{G}$  is expressed in terms of the functions  $\hat{\varphi}_i$ , as in [8, Eq. (2.6)]. Moreover, the functions  $\hat{\varphi}_i$  themselves are usually tensor-product B-splines or NURBS. In this paper, the role of the  $\hat{\varphi}_i$  will be played by tensor-product B-splines over uniform knot sequences. Furthermore, we do not confine ourselves to the isoparametric approach, but we also allow the geometry map  $\mathbf{G}$  to be any sufficiently regular function from  $[0, 1]^d$  to  $\overline{\Omega}$ , not necessarily expressed in terms of B-splines.

### 3.2 B-spline basis functions and IgA Galerkin matrices.

Let us now provide the explicit construction of our basis functions  $\hat{\varphi}_i$ . For  $p, n \geq 1$ , consider the uniform knot sequence

$$t_1 = \dots = t_{p+1} = 0 < t_{p+2} < \dots < t_{p+n} < 1 = t_{p+n+1} = \dots = t_{2p+n+1},$$

where

$$t_{i+p+1} := \frac{i}{n}, \quad i = 0, \dots, n.$$

The B-splines of degree  $p$  on this knot sequence are denoted by

$$N_{i,[p]} : [0, 1] \rightarrow \mathbb{R}, \quad i = 1, \dots, n + p,$$

and are defined recursively as follows [3]: for  $1 \leq i \leq (n + p) + p$ ,

$$N_{i,[0]}(x) := \begin{cases} 1, & \text{if } x \in [t_i, t_{i+1}), \\ 0, & \text{elsewhere;} \end{cases}$$

for  $1 \leq k \leq p$  and  $1 \leq i \leq (n + p) + p - k$ ,

$$N_{i,[k]}(x) := \frac{x - t_i}{t_{i+k} - t_i} N_{i,[k-1]}(x) + \frac{t_{i+k+1} - x}{t_{i+k+1} - t_{i+1}} N_{i+1,[k-1]}(x),$$

where we assume that a fraction with zero denominator is zero. We know from [3] that the functions  $N_{1,[p]}, \dots, N_{n+p,[p]}$  form a basis for the spline space

$$\{s \in C^{p-1}([0, 1]) : s|_{[\frac{i}{n}, \frac{i+1}{n}]} \in \mathbb{P}_p, \quad \forall i = 0, \dots, n-1\},$$

where  $\mathbb{P}_p$  is the space of polynomials of degree less than or equal to  $p$ . Moreover, the functions  $N_{i,[p]}$  have a local support, namely

$$\text{supp}(N_{i,[p]}) = [t_i, t_{i+p+1}], \quad i = 1, \dots, n + p,$$

and

$$N_{i,[p]}(0) = N_{i,[p]}(1) = 0, \quad i = 2, \dots, n + p - 1.$$

Now, let  $\mathbf{p} := (p_1, \dots, p_d)$  and  $\mathbf{n} := (n_1, \dots, n_d)$  be multi-indices in  $\mathbb{N}^d$ , and define the tensor-product B-splines  $N_{i,[\mathbf{p}]} : [0, 1]^d \rightarrow \mathbb{R}$  by

$$N_{i,[\mathbf{p}]} := N_{i_1,[p_1]} \otimes \dots \otimes N_{i_d,[p_d]}, \quad \mathbf{i} = \mathbf{2}, \dots, \mathbf{n} + \mathbf{p} - \mathbf{1}. \quad (3.6)$$

In the framework of IgA based on (uniform) B-splines, the functions  $\hat{\varphi}_i$ ,  $i = 1, \dots, N$ , in (3.2) are chosen as the tensor-product B-splines in (3.6). In this setting,  $N = N(\mathbf{n} + \mathbf{p} - \mathbf{2})$ . Moreover, we adopt for the tensor-product B-splines (3.6) the standard lexicographic ordering (2.1).<sup>1</sup> This ordering is followed when assembling the stiffness matrix (3.4), which from now on will be denoted by  $A_{\mathbf{G},\mathbf{n}}^{[\mathbf{p}]}$ , in order to emphasize its dependence on  $\mathbf{p}, \mathbf{n}$  and the geometry map  $\mathbf{G}$ . In multi-index notation, the matrix  $A_{\mathbf{G},\mathbf{n}}^{[\mathbf{p}]}$  is expressed by

$$A_{\mathbf{G},\mathbf{n}}^{[\mathbf{p}]} := \left[ \int_{[0,1]^d} ((\nabla N_{j+1,[\mathbf{p}]})^T K_{\mathbf{G}} \nabla N_{i+1,[\mathbf{p}]} + (\nabla N_{j+1,[\mathbf{p}]})^T (J_{\mathbf{G}})^{-1} \boldsymbol{\alpha}(\mathbf{G}) N_{i+1,[\mathbf{p}]} + \gamma(\mathbf{G}) N_{j+1,[\mathbf{p}]} N_{i+1,[\mathbf{p}]}) | \det(J_{\mathbf{G}}) | \right]_{i,j=1}^{n+p-2}.$$

Note that  $A_{\mathbf{G},\mathbf{n}}^{[\mathbf{p}]}$  can be decomposed as follows:

$$A_{\mathbf{G},\mathbf{n}}^{[\mathbf{p}]} = K_{\mathbf{G},\mathbf{n}}^{[\mathbf{p}]} + R_{\mathbf{G},\mathbf{n}}^{[\mathbf{p}]}, \quad (3.7)$$

<sup>1</sup>Although the lexicographic ordering is very common in the literature, an alternative ordering has been used in [6, 7, 8, 11].

where

$$K_{\mathbf{G},\mathbf{n}}^{[p]} := \left[ \int_{[0,1]^d} (\nabla N_{j+1,[p]})^T K_{\mathbf{G}} |\det(J_{\mathbf{G}})| \nabla N_{i+1,[p]} \right]_{i,j=1}^{n+p-2} \quad (3.8)$$

is the matrix resulting from the discretization of the diffusive term in (1.1), and

$$R_{\mathbf{G},\mathbf{n}}^{[p]} := \left[ \int_{[0,1]^d} ((\nabla N_{j+1,[p]})^T (J_{\mathbf{G}})^{-1} \alpha(\mathbf{G}) N_{i+1,[p]} + \gamma(\mathbf{G}) N_{j+1,[p]} N_{i+1,[p]}) |\det(J_{\mathbf{G}})| \right]_{i,j=1}^{n+p-2} \quad (3.9)$$

is the matrix resulting from the discretization of the terms in (1.1) with lower order derivatives. The matrix  $R_{\mathbf{G},\mathbf{n}}^{[p]}$  can be regarded as a ‘residual term’. Indeed, the norm of  $R_{\mathbf{G},\mathbf{n}}^{[p]}$  is negligible with respect to the norm of the diffusion matrix  $K_{\mathbf{G},\mathbf{n}}^{[p]}$  when the discretization parameters  $\mathbf{n}$  are large; see Section 4.2 (Step 1) and Section 4.3.

## 4 Spectral distribution

### 4.1 The symbol of the normalized IgA Galerkin matrices

Let  $\mathbb{Q}_+^d := \{\mathbf{r} := (r_1, \dots, r_d) \in \mathbb{Q}^d : r_i > 0, \forall i = 1, \dots, d\}$  and fix a vector  $\mathbf{v} := (v_1, \dots, v_d) \in \mathbb{Q}_+^d$ . From now on, until the end of the paper, we assume that  $n_j = v_j n$  for every  $j = 1, \dots, d$ , i.e.,  $\mathbf{n} = \mathbf{v}n$ . It is understood that  $n$  varies in the set of indices such that  $\mathbf{v}n \in \mathbb{N}^d$ .

In our main result (Theorem 4.1), we consider the sequence of normalized matrices  $\{n^{d-2} A_{\mathbf{G},\mathbf{n}}^{[p]}\}_n$ , and we compute its singular value and eigenvalue distribution in the sense of Definition 2.1. The proof of Theorem 4.1 heavily relies on the theory of sLT and GLT sequences developed in [20, 21]. Actually, we show that  $\{n^{d-2} A_{\mathbf{G},\mathbf{n}}^{[p]}\}_n$  can be expressed as a finite sum of sLT sequences and is therefore a GLT sequence. Despite the importance of sLT and GLT sequences, we do not present the corresponding general definitions, since they are rather difficult, and we refer the reader to [20, Definition 2.1] for the definition of sLT sequences and to [20, Definition 2.3] for the definition of GLT sequences. As we shall see in the following, we can somehow limit the use of the formal definitions to some basic cases and prove Theorem 4.1 by exploiting directly the properties of sLT and GLT sequences provided in [20, 21].

Let us start with defining the following  $d \times d$  symmetric matrix  $H_p$ , whose components are continuous functions in the Fourier variables  $\boldsymbol{\theta} := (\theta_1, \dots, \theta_d) \in [-\pi, \pi]^d$ :

$$(H_p)_{ij} := \begin{cases} \left( \bigotimes_{r=1}^{i-1} h_{p_r} \right) \otimes f_{p_i} \otimes \left( \bigotimes_{r=i+1}^d h_{p_r} \right), & \text{if } i = j, \\ \left( \bigotimes_{r=1}^{i-1} h_{p_r} \right) \otimes g_{p_i} \otimes \left( \bigotimes_{r=i+1}^{j-1} h_{p_r} \right) \otimes g_{p_j} \otimes \left( \bigotimes_{r=j+1}^d h_{p_r} \right), & \text{if } i < j, \\ \left( \bigotimes_{r=1}^{j-1} h_{p_r} \right) \otimes g_{p_j} \otimes \left( \bigotimes_{r=j+1}^{i-1} h_{p_r} \right) \otimes g_{p_i} \otimes \left( \bigotimes_{r=i+1}^d h_{p_r} \right), & \text{if } i > j, \end{cases} \quad (4.1)$$

where  $h_p, g_p, f_p : [-\pi, \pi] \rightarrow \mathbb{R}$  are defined for all  $p \geq 1$  by

$$h_p(\theta) := \phi_{[2p+1]}(p+1) + 2 \sum_{k=1}^p \phi_{[2p+1]}(p+1-k) \cos(k\theta), \quad (4.2)$$

$$g_p(\theta) := -2 \sum_{k=1}^p \dot{\phi}_{[2p+1]}(p+1-k) \sin(k\theta), \quad (4.3)$$

$$f_p(\theta) := -\ddot{\phi}_{[2p+1]}(p+1) - 2 \sum_{k=1}^p \ddot{\phi}_{[2p+1]}(p+1-k) \cos(k\theta), \quad (4.4)$$

and  $\phi_{[p]}, \dot{\phi}_{[p]}, \ddot{\phi}_{[p]}$  are, respectively, the cardinal B-spline of degree  $p$  (defined, e.g., in [11, Eqs. (7)–(8)]), its first derivative and its second derivative. The  $(i, j)$ -th entry of the matrix  $H_p$  is related to the second order partial derivative  $-\frac{\partial^2}{\partial \hat{x}_i \partial \hat{x}_j}$ .

**Theorem 4.1.** *Assume that the geometry map  $\mathbf{G}$  is regular, i.e.,  $\mathbf{G} \in C^1([0, 1]^d)$  and  $\det(J_{\mathbf{G}}) \neq 0$  in  $[0, 1]^d$ , and suppose that the components of  $K$  are continuous over  $\bar{\Omega}$ . Then, the sequence of normalized matrices  $\{n^{d-2}A_{\mathbf{G},n}^{[p]}\}_n$  is distributed, in the sense of both the singular values and the eigenvalues, like the function  $f_{\mathbf{G},p}^{(\nu)} : [0, 1]^d \times [-\pi, \pi]^d \rightarrow \mathbb{R}$ ,*

$$f_{\mathbf{G},p}^{(\nu)}(\hat{\mathbf{x}}, \boldsymbol{\theta}) := \frac{\boldsymbol{\nu} \left( |\det(J_{\mathbf{G}}(\hat{\mathbf{x}}))| K_{\mathbf{G}}(\hat{\mathbf{x}}) \circ H_p(\boldsymbol{\theta}) \right) \boldsymbol{\nu}^T}{N(\boldsymbol{\nu})}, \quad (4.5)$$

where  $\circ$  is the componentwise (Hadamard) product of matrices and  $H_p$  is defined in (4.1). In formulas,  $\{n^{d-2}A_{\mathbf{G},n}^{[p]}\}_n \sim_{\sigma} f_{\mathbf{G},p}^{(\nu)}$  and  $\{n^{d-2}A_{\mathbf{G},n}^{[p]}\}_n \sim_{\lambda} f_{\mathbf{G},p}^{(\nu)}$ .

Before going into the details of the proof of Theorem 4.1, we first outline the main idea in Section 4.2. This idea is quite general and provides an abstract framework for the computation of the singular value and eigenvalue distribution of matrix-sequences coming from a PDE discretization. It can be also applied to other approximation techniques besides Galerkin B-spline IgA. For example, we can refer to FDs [20, Section 6], FEs [1, 13], and isogeometric collocation methods [8].

*Remark 4.2.* The proof of Theorem 4.1 does not require that  $K$  is positive definite. Only the symmetry of  $K$  is needed. However, the positive definiteness of the diffusion matrix  $K$  is necessary for the ellipticity of the differential problem (1.1).

*Remark 4.3.* In order to simplify the presentation of the proof of Theorem 4.1, the components of  $K$  are required to be continuous. Actually, the continuity requirement can be relaxed by assuming that the components of  $K$  are only piecewise continuous. With a little more effort, one could also extend the validity of Theorem 4.1 to the case where the components of  $K(\mathbf{G})$  are Riemann integrable, taking into consideration that Riemann integrability is the classical assumption under which the singular value and eigenvalue distribution results for sLT and GLT sequences have been formulated; see, e.g., [20, Theorems 4.1–4.8].

*Remark 4.4.* The real challenge is to show that Theorem 4.1 holds even if the components of  $K$  are only in  $L^\infty(\Omega)$ , as asserted in the formulation of our problem (1.1). The extension of the proof to the  $L^\infty$  case would require the application of the Lusin theorem [17], which is used in function theory for approximating a measurable function by a continuous function. We refer the reader to [19, Theorem 6.3] for an application of the Lusin theorem in the context of FDs under homogenization. The extension of Theorem 4.1 to the  $L^\infty$  case is an interesting topic for further investigation.

## 4.2 The idea for proving Theorem 4.1

The idea consists of the following steps, which will be detailed afterwards. Here, we just outline the key ingredients, so as to make more clear the generality of the argument we are going to see.

We first recall that, for any two sequences of real numbers  $\{a_n\}_n$  and  $\{b_n\}_n$ , the notation  $a_n = O(b_n)$  means that there exists a constant  $C$ , independent of  $n$ , such that  $|a_n| \leq C|b_n|$  for all  $n$ . On the other hand,  $a_n = o(b_n)$  means  $\lim_{n \rightarrow \infty} a_n/b_n = 0$ . A sequence of real numbers  $\{a_n\}_n$  is said to be bounded away from 0 (resp.  $\infty$ ) if there exists a positive constant  $c$  such that  $|a_n| \geq c$  (resp.  $|a_n| \leq c$ ) for all  $n$ . Note that the terminology ‘bounded away from  $\infty$ ’ is equivalent to ‘uniformly bounded with respect to  $n$ ’.

**Step 1.** By definition, see (3.9),  $R_{\mathbf{G},n}^{[p]}$  is the matrix resulting from the discretization of the terms in (1.1) with lower order derivatives. Since  $n^{d-2}$  is the ‘correct’ normalization factor, which keeps the spectral norm of the diffusion matrix  $n^{d-2}K_{\mathbf{G},n}^{[p]}$  bounded away from 0 and  $\infty$ , we can show that  $\|n^{d-2}R_{\mathbf{G},n}^{[p]}\| \rightarrow 0$  as  $n \rightarrow \infty$ , and more precisely  $\|n^{d-2}R_{\mathbf{G},n}^{[p]}\| = O(n^{-1})$ .

By (2.3) and by the definition of sLT sequences [20, Definition 2.1], any sequence of matrices  $\{X_n\}_n$  with increasing dimension such that  $\|X_n\| \rightarrow 0$  is an sLT sequence with the so-called weight function and generating function both equal to 0; this is denoted by  $\{X_n\}_n \sim_{\text{sLT}} (0, 0)$ . Hence, we have  $\{n^{d-2}R_{\mathbf{G},n}^{[p]}\}_n \sim_{\text{sLT}} (0, 0)$ .

$(0, 0)$ . Moreover, from the definition of GLT sequences [20, Definition 2.3] it follows immediately that the relation  $\{X_n\}_n \sim_{\text{sLT}} (a, f)$ , with  $a : [0, 1]^d \rightarrow \mathbb{C}$  (measurable) and  $f : [-\pi, \pi]^d \rightarrow \mathbb{C}$ , implies  $\{X_n\}_n \sim_{\text{GLT}} a \otimes f = a(\hat{\mathbf{x}})f(\boldsymbol{\theta})$ . The function  $a \otimes f$  is referred to as the symbol of the GLT sequence  $\{X_n\}_n$ . Thus, we have  $\{n^{d-2}K_{\mathbf{G},n}^{[p]}\}_n \sim_{\text{GLT}} 0$ .

Finally, the GLT sequences form an algebra and, in particular, any linear combination of GLT sequences is still a GLT sequence, whose symbol is given by the same linear combination of the involved symbols; see [21, Theorem 2.2 (especially the first two lines of the proof)]. Hence, by (3.7),  $\{n^{d-2}A_{\mathbf{G},n}^{[p]}\}_n \sim_{\text{GLT}} f_{\mathbf{G},p}^{(\nu)}$  if and only if  $\{n^{d-2}K_{\mathbf{G},n}^{[p]}\}_n \sim_{\text{GLT}} f_{\mathbf{G},p}^{(\nu)}$ . In this way, we have reduced the analysis of the sequence  $\{n^{d-2}A_{\mathbf{G},n}^{[p]}\}_n$  to the analysis of its strictly diffusive part  $\{n^{d-2}K_{\mathbf{G},n}^{[p]}\}_n$ .

**Step 2.** Let  $[L^1([0, 1]^d)]^{d \times d}$  be the space of functions  $L : [0, 1]^d \rightarrow \mathbb{R}^{d \times d}$  such that  $L_{ij} \in L^1([0, 1]^d)$  for all  $i, j = 1, \dots, d$ . Let us consider the operator  $\mathcal{L}_n^{[p]}(\cdot) : [L^1([0, 1]^d)]^{d \times d} \rightarrow \mathbb{R}^{N(n+p-2) \times N(n+p-2)}$ ,

$$\mathcal{L}_n^{[p]}(L) := \left[ \int_{[0,1]^d} (\nabla N_{j+1,[p]})^T L \nabla N_{i+1,[p]} \right]_{i,j=1}^{n+p-2}. \quad (4.6)$$

By (3.8), we have

$$K_{\mathbf{G},n}^{[p]} = \mathcal{L}_n^{[p]}(K_{\mathbf{G}} | \det(J_{\mathbf{G}})|).$$

During the next steps, we shall use the linearity of  $\mathcal{L}_n^{[p]}(\cdot)$  as well as the algebra structure of GLT sequences (i.e., its closure under linear combinations and products) to prove that  $\{n^{d-2}K_{\mathbf{G},n}^{[p]}\}_n \sim_{\text{GLT}} f_{\mathbf{G},p}^{(\nu)}$ .

**Step 3.** Take  $L = E_{st}$  in (4.6), where  $E_{st}$  is the  $d \times d$  matrix having 1 in position  $(s, t)$  and 0 elsewhere. Note that if  $\mathbf{G}$  is the identity map and we take  $K = E_{st}$  in (1.1), then we are ‘selecting’ the second order partial derivative  $-\frac{\partial^2 u}{\partial x_s \partial x_t} = -\frac{\partial^2 u}{\partial \hat{x}_s \partial \hat{x}_t}$ , which is a separable differential operator.<sup>2</sup> With a direct computation, one can show that  $n^{d-2} \mathcal{L}_n^{[p]}(E_{st})$  is the  $d$ -level Toeplitz matrix

$$\frac{\nu_s \nu_t}{N(\boldsymbol{\nu})} T_{n+p-2}((H_p)_{st}), \quad (4.7)$$

up to a correction whose rank is  $O(n^{d-1}) = O(N(n+p-2)/n)$ . In the expression (4.7),  $H_p(\boldsymbol{\theta})$  is the  $d \times d$  symmetric matrix in the Fourier variables  $\boldsymbol{\theta} \in [-\pi, \pi]^d$ , whose  $(i, j)$ -th entry comes from the discretization of the second order partial derivative  $-\frac{\partial^2}{\partial \hat{x}_i \partial \hat{x}_j}$ . In our specific case, this matrix is given by (4.1).

Therefore, using the definition of sLT sequences one can show that

$$\{n^{d-2} \mathcal{L}_n^{[p]}(E_{st})\}_n \sim_{\text{sLT}} \left( 1, \frac{\nu_s \nu_t}{N(\boldsymbol{\nu})} (H_p)_{st} \right),$$

and so

$$\{n^{d-2} \mathcal{L}_n^{[p]}(E_{st})\}_n \sim_{\text{GLT}} 1 \otimes \frac{\nu_s \nu_t}{N(\boldsymbol{\nu})} (H_p)_{st} = \frac{\nu_s \nu_t}{N(\boldsymbol{\nu})} (H_p)_{st}.$$

<sup>2</sup>We say that a differential operator is separable if it is obtained by multiplying a given function with a product of partial derivatives. The general separable differential operator can be written as

$$a(\mathbf{x}) \frac{\partial^{r_1 + \dots + r_d} u}{\partial x_1^{r_1} \dots \partial x_d^{r_d}}.$$

An example of a non-separable differential operator is the Laplacian, which however can be written, like any other differential operator, as a sum of separable differential operators:

$$\Delta u = \sum_{k=1}^d \frac{\partial^2 u}{\partial x_k^2}.$$

What we are going to see is that a separable differential operator gives rise to an sLT sequence. As a consequence, an arbitrary differential operator (a sum of separable differential operators) gives rise to a sum of sLT sequences, i.e., a GLT sequence.

If  $L$  is a constant matrix, then  $L = \sum_{s,t=1}^d L_{st} E_{st}$ , and by the closure of GLT sequences with respect to linear combinations we have

$$\{n^{d-2} \mathcal{L}_n^{[p]}(L)\}_n \sim_{\text{GLT}} \sum_{s,t=1}^d L_{st} \frac{v_s v_t}{N(\mathbf{v})} (H_p)_{st} = \frac{\mathbf{v}(L \circ H_p) \mathbf{v}^T}{N(\mathbf{v})}.$$

**Step 4.** Let us pass to the variable coefficient case and, inspired by Step 3, let us take  $L = a(\hat{\mathbf{x}}) E_{st}$  in (4.6), where  $a$  is a continuous function defined on  $[0, 1]^d$ . Note that if  $\mathbf{G}$  is the identity map and we take  $K = a(\hat{\mathbf{x}}) E_{st}$  in (1.1), then we are ‘selecting’  $-\frac{\partial}{\partial \hat{x}_s} (a(\hat{\mathbf{x}}) \frac{\partial u}{\partial \hat{x}_t})$ , which coincides with the second order separable differential operator  $-a(\hat{\mathbf{x}}) \frac{\partial^2 u}{\partial \hat{x}_s \partial \hat{x}_t}$  up to a term with a lower order derivative, namely  $-\frac{\partial a}{\partial \hat{x}_s}(\hat{\mathbf{x}}) \frac{\partial u}{\partial \hat{x}_t}$ . Given the local support of the basis functions  $N_{i, [p]}$ ,  $i = 2, \dots, n+p-1$ , and the fact that  $\text{supp}(N_{i, [p]})$  is located near the point  $\mathbf{i}/\mathbf{n} = (i_1/n_1, \dots, i_d/n_d)$ , it can be shown that

$$n^{d-2} \mathcal{L}_n^{[p]}(a(\hat{\mathbf{x}}) E_{st}) = n^{d-2} D_{n+p-2}(a) \mathcal{L}_n^{[p]}(E_{st}) + Q_{n+p-2}. \quad (4.8)$$

The matrix  $Q_{n+p-2}$  is a (sparse) matrix whose components are  $O(\omega_a(n^{-1}))$  as well as its spectral norm. Here,  $\omega_a(\cdot)$  stands for the modulus of continuity of the function  $a$ . Using the decomposition (4.8), we get

$$\{n^{d-2} \mathcal{L}_n^{[p]}(a(\hat{\mathbf{x}}) E_{st})\}_n \sim_{\text{sLT}} \left( a, \frac{v_s v_t}{N(\mathbf{v})} (H_p)_{st} \right),$$

and so

$$\{n^{d-2} \mathcal{L}_n^{[p]}(a(\hat{\mathbf{x}}) E_{st})\}_n \sim_{\text{GLT}} a \otimes \frac{v_s v_t}{N(\mathbf{v})} (H_p)_{st} = \frac{v_s v_t}{N(\mathbf{v})} a(\hat{\mathbf{x}}) (H_p(\boldsymbol{\theta}))_{st}.$$

**Step 5.** To obtain the relation  $\{n^{d-2} K_{\mathbf{G}, \mathbf{n}}^{[p]}\}_n \sim_{\text{GLT}} f_{\mathbf{G}, \mathbf{p}}^{(\mathbf{v})}$ , we invoke the linearity of  $\mathcal{L}_n^{[p]}(\cdot)$  and the closure of the GLT sequences under linear combinations. Since

$$K_{\mathbf{G}} |\det(J_{\mathbf{G}})| = \sum_{s,t=1}^d (K_{\mathbf{G}} |\det(J_{\mathbf{G}})|)_{st} E_{st},$$

with  $(K_{\mathbf{G}} |\det(J_{\mathbf{G}})|)_{st} \in C([0, 1]^d)$  for every  $s, t = 1, \dots, d$ , we have

$$n^{d-2} \mathcal{L}_n^{[p]}(K_{\mathbf{G}} |\det(J_{\mathbf{G}})|) = \sum_{s,t=1}^d n^{d-2} \mathcal{L}_n^{[p]}((K_{\mathbf{G}} |\det(J_{\mathbf{G}})|)_{st} E_{st}). \quad (4.9)$$

This shows that  $\{n^{d-2} K_{\mathbf{G}, \mathbf{n}}^{[p]}\}_n = \{n^{d-2} \mathcal{L}_n^{[p]}(K_{\mathbf{G}} |\det(J_{\mathbf{G}})|)\}_n$  is a finite sum of sLT sequences. Moreover, using the closure of GLT sequences under linear combinations, we get

$$\begin{aligned} \{n^{d-2} \mathcal{L}_n^{[p]}(K_{\mathbf{G}} |\det(J_{\mathbf{G}})|)\}_n &\sim_{\text{GLT}} \sum_{s,t=1}^d \frac{v_s v_t}{N(\mathbf{v})} (K_{\mathbf{G}}(\hat{\mathbf{x}}) |\det(J_{\mathbf{G}}(\hat{\mathbf{x}}))|)_{st} (H_p(\boldsymbol{\theta}))_{st} \\ &= \frac{\mathbf{v} (K_{\mathbf{G}}(\hat{\mathbf{x}}) |\det(J_{\mathbf{G}}(\hat{\mathbf{x}}))| \circ H_p(\boldsymbol{\theta})) \mathbf{v}^T}{N(\mathbf{v})}, \end{aligned}$$

which is nothing else than  $\{n^{d-2} K_{\mathbf{G}, \mathbf{n}}^{[p]}\}_n \sim_{\text{GLT}} f_{\mathbf{G}, \mathbf{p}}^{(\mathbf{v})}$ .

**Step 6.** The singular value distribution  $\{n^{d-2} A_{\mathbf{G}, \mathbf{n}}^{[p]}\}_n \sim_{\sigma} f_{\mathbf{G}, \mathbf{p}}^{(\mathbf{v})}$  follows from the relation  $\{n^{d-2} A_{\mathbf{G}, \mathbf{n}}^{[p]}\}_n \sim_{\text{GLT}} f_{\mathbf{G}, \mathbf{p}}^{(\mathbf{v})}$  and from either [20, Theorem 4.5] or [20, Theorem 4.1]. If the matrices  $n^{d-2} A_{\mathbf{G}, \mathbf{n}}^{[p]}$  are symmetric (this happens when  $\boldsymbol{\alpha} = \mathbf{0}$ ), then the eigenvalue distribution  $\{n^{d-2} A_{\mathbf{G}, \mathbf{n}}^{[p]}\}_n \sim_{\lambda} f_{\mathbf{G}, \mathbf{p}}^{(\mathbf{v})}$  follows from the relation  $\{n^{d-2} A_{\mathbf{G}, \mathbf{n}}^{[p]}\}_n \sim_{\text{GLT}}$

$f_{\mathbf{G},\mathbf{p}}^{(\nu)}$  and from either [20, Theorem 4.8] or [20, Theorem 4.4]. When applying all these theorems, take also into account that  $\{n^{d-2}A_{\mathbf{G},n}^{[p]}\}_n$  is the sum of  $d^2 + 1$  sLT sequences with continuous (and hence Riemann integrable) weight functions: the  $d^2$  sLT sequences in (4.9), whose sum is  $\{n^{d-2}K_{\mathbf{G},n}^{[p]}\}_n$ , plus the sLT sequence  $\{n^{d-2}R_{\mathbf{G},n}^{[p]}\}_n \sim_{\text{sLT}} (0, 0)$ . In the case where the matrices  $n^{d-2}A_{\mathbf{G},n}^{[p]}$  are not symmetric, the eigenvalue distribution  $\{n^{d-2}A_{\mathbf{G},n}^{[p]}\}_n \sim_{\lambda} f_{\mathbf{G},\mathbf{p}}^{(\nu)}$  still holds. Indeed,

- a)  $\{n^{d-2}K_{\mathbf{G},n}^{[p]}\}_n \sim_{\lambda} f_{\mathbf{G},\mathbf{p}}^{(\nu)}$ , because  $\{n^{d-2}K_{\mathbf{G},n}^{[p]}\}_n \sim_{\text{GLT}} f_{\mathbf{G},\mathbf{p}}^{(\nu)}$  and the matrices  $K_{\mathbf{G},n}^{[p]}$  are symmetric;
- b)  $\|n^{d-2}K_{\mathbf{G},n}^{[p]}\|$  is uniformly bounded with respect to  $n$ , due to the correct normalization factor  $n^{d-2}$  (see Step 1);
- c)  $\|n^{d-2}R_{\mathbf{G},n}^{[p]}\| \rightarrow 0$  (see Step 1) and, consequently,  $\|n^{d-2}R_{\mathbf{G},n}^{[p]}\|_1 = o(N(\mathbf{n} + \mathbf{p} - \mathbf{2}))$  by (2.3).

Therefore, all the assumptions of [15, Theorem 3.4] or of its generalized version [12, Theorem 3.3] are satisfied, and the eigenvalue distribution  $\{n^{d-2}A_{\mathbf{G},n}^{[p]}\}_n \sim_{\lambda} f_{\mathbf{G},\mathbf{p}}^{(\nu)}$  follows.

### 4.3 Proof of Step 1

We show that  $\|n^{d-2}R_{\mathbf{G},n}^{[p]}\| = O(n^{-1})$ . We also prove that  $\|n^{d-2}K_{\mathbf{G},n}^{[p]}\|$  is uniformly bounded with respect to  $n$ : this is asserted in Step 1 and is used in Step 6, item (b). In the following,  $C$  denotes a generic constant independent of  $n$ .

By hypothesis,  $\mathbf{G}$  is regular, so the components of  $(J_{\mathbf{G}})^{-1}$  are continuous and bounded over  $[0, 1]^d$ . Moreover, the coefficients  $\gamma$  and  $\alpha_i$ ,  $i = 1, \dots, d$ , in (1.1) are bounded. Hence, by (3.9), for all  $\mathbf{i}, \mathbf{j} = \mathbf{2}, \dots, \mathbf{n} + \mathbf{p} - \mathbf{1}$  we have

$$\begin{aligned} |(R_{\mathbf{G},n}^{[p]})_{\mathbf{i}-1, \mathbf{j}-1}| &\leq \int_{[0,1]^d} \left( |(\nabla N_{\mathbf{j},[\mathbf{p}]})^T (J_{\mathbf{G}})^{-1} \alpha(\mathbf{G}) N_{\mathbf{i},[\mathbf{p}]}| + |\gamma(\mathbf{G}) N_{\mathbf{j},[\mathbf{p}]} N_{\mathbf{i},[\mathbf{p}]}| \right) |\det(J_{\mathbf{G}})| \\ &\leq C \left[ \int_{[0,1]^d} \sum_{r=1}^d \left| \frac{\partial N_{\mathbf{j},[\mathbf{p}]}}{\partial \hat{x}_r} \right| |N_{\mathbf{i},[\mathbf{p}]}| + \int_{[0,1]^d} |N_{\mathbf{j},[\mathbf{p}]}| |N_{\mathbf{i},[\mathbf{p}]}| \right] \\ &= C \left[ \sum_{r=1}^d \int_0^1 |N'_{j_r, [p_r]}| |N_{i_r, [p_r]}| \prod_{\substack{s=1 \\ s \neq r}}^d \int_0^1 |N_{j_s, [p_s]}| |N_{i_s, [p_s]}| + \prod_{s=1}^d \int_0^1 |N_{j_s, [p_s]}| |N_{i_s, [p_s]}| \right]. \end{aligned}$$

From the positivity of the B-splines  $N_{i, [\mathbf{p}]}$ , the results in [11, proof of Lemma 8] and the relation  $\mathbf{n} = \nu \mathbf{n}$ , we get

$$|(R_{\mathbf{G},n}^{[p]})_{\mathbf{i}-1, \mathbf{j}-1}| \leq C \left[ \sum_{r=1}^d 2 \prod_{\substack{s=1 \\ s \neq r}}^d \frac{1}{n_s} + \prod_{s=1}^d \frac{1}{n_s} \right] = O(n^{-d+1}).$$

In view of the local support property of the B-splines,  $\text{supp}(N_{i, [\mathbf{p}]}) = [t_i, t_{i+p+1}]$ , it is clear that  $(R_{\mathbf{G},n}^{[p]})_{\mathbf{i}-1, \mathbf{j}-1} = 0$  if  $\|\mathbf{i} - \mathbf{j}\|_{\infty} \geq \|\mathbf{p}\|_{\infty} + 1$ . It follows that  $R_{\mathbf{G},n}^{[p]}$  is a sparse matrix, consisting of at most  $(2\|\mathbf{p}\|_{\infty} + 1)^d$  (independent of  $n$ ) non-zero components in each row and column with a value  $O(n^{-d+1})$ . Thus, by (2.2), we have  $\|R_{\mathbf{G},n}^{[p]}\| = O(n^{-d+1})$  and  $\|n^{d-2}R_{\mathbf{G},n}^{[p]}\| = O(n^{-1})$ . In particular,  $\|n^{d-2}R_{\mathbf{G},n}^{[p]}\| \rightarrow 0$  when  $n \rightarrow \infty$ .

To prove that  $\|n^{d-2}K_{\mathbf{G},n}^{[p]}\|$  is uniformly bounded with respect to  $n$ , we can follow the same pattern as for

the proof of  $\|n^{d-2}R_{G,n}^{[p]}\| = O(n^{-1})$ . By (3.8), for all  $i, j = 2, \dots, n+p-1$  we have

$$\begin{aligned}
|(K_{G,n}^{[p]})_{i-1,j-1}| &\leq \int_{[0,1]^d} |(\nabla N_{j,[p]})^T K_G \nabla N_{i,[p]}| |\det(J_G)| \\
&\leq C \int_{[0,1]^d} \sum_{r,s=1}^d \left| \frac{\partial N_{j,[p]}}{\partial \hat{x}_r} \right| \left| \frac{\partial N_{i,[p]}}{\partial \hat{x}_s} \right| \\
&= C \left[ \sum_{r=1}^d \int_0^1 |N'_{j_r,[p_r]}| |N'_{i_r,[p_r]}| \prod_{\substack{t=1 \\ t \neq r}}^d \int_0^1 |N_{j_t,[p_t]}| |N_{i_t,[p_t]}| \right. \\
&\quad \left. + \sum_{\substack{r,s=1 \\ r \neq s}}^d \int_0^1 |N'_{j_r,[p_r]}| |N_{i_r,[p_r]}| \int_0^1 |N_{j_s,[p_s]}| |N'_{i_s,[p_s]}| \prod_{\substack{t=1 \\ t \neq r,s}}^d \int_0^1 |N_{j_t,[p_t]}| |N_{i_t,[p_t]}| \right] \\
&\leq C \left[ \sum_{r=1}^d 4p_r n_r \prod_{\substack{t=1 \\ t \neq r}}^d \frac{1}{n_t} + \sum_{\substack{r,s=1 \\ r \neq s}}^d 2 \cdot 2 \prod_{\substack{t=1 \\ t \neq r,s}}^d \frac{1}{n_t} \right] = O(n^{-d+2}),
\end{aligned}$$

where in the last inequality we used again the positivity of the B-splines  $N_{i,[p]}$  and the results in [11, proof of Lemma 8]. By means of the local support of the basis functions  $N_{i,[p]}$  and by following the same argument used for  $R_{G,n}^{[p]}$ , it can be shown that  $K_{G,n}^{[p]}$  is a sparse matrix, with at most  $(2\|\mathbf{p}\|_\infty + 1)^d$  non-zero components in each row and column. Thus,  $\|n^{d-2}K_{G,n}^{[p]}\| = O(1)$  is uniformly bounded with respect to  $n$ .

Finally, we note that in Step 2 there is nothing to prove.

#### 4.4 Proof of Step 3

Let  $1 \leq s, t \leq d$ . We first show that the rank of the difference

$$n^{d-2} \mathcal{L}_n^{[p]}(E_{st}) - \frac{\nu_s \nu_t}{N(\mathbf{v})} T_{n+p-2}((H_p)_{st}) \quad (4.10)$$

is  $O(n^{d-1})$  as  $n \rightarrow \infty$ , with  $H_p$  defined in (4.1). Then, we prove that

$$\{n^{d-2} \mathcal{L}_n^{[p]}(E_{st})\}_n \sim_{\text{sLT}} \left( 1, \frac{\nu_s \nu_t}{N(\mathbf{v})} (H_p)_{st} \right). \quad (4.11)$$

By following the construction of the matrices in [10, Eqs. (3.17), (4.38) and (5.17)], one can show that

$$\mathcal{L}_n^{[p]}(E_{st}) = \begin{cases} \left( \bigotimes_{r=1}^{s-1} \frac{1}{n_r} M_{n_r}^{[p_r]} \right) \otimes n_s K_{n_s}^{[p_s]} \otimes \left( \bigotimes_{r=s+1}^d \frac{1}{n_r} M_{n_r}^{[p_r]} \right), & \text{if } s = t, \\ - \left( \bigotimes_{r=1}^{s-1} \frac{1}{n_r} M_{n_r}^{[p_r]} \right) \otimes H_{n_s}^{[p_s]} \otimes \left( \bigotimes_{r=s+1}^{t-1} \frac{1}{n_r} M_{n_r}^{[p_r]} \right) \otimes H_{n_t}^{[p_t]} \otimes \left( \bigotimes_{r=t+1}^d \frac{1}{n_r} M_{n_r}^{[p_r]} \right), & \text{if } s < t, \\ - \left( \bigotimes_{r=1}^{t-1} \frac{1}{n_r} M_{n_r}^{[p_r]} \right) \otimes H_{n_t}^{[p_t]} \otimes \left( \bigotimes_{r=t+1}^{s-1} \frac{1}{n_r} M_{n_r}^{[p_r]} \right) \otimes H_{n_s}^{[p_s]} \otimes \left( \bigotimes_{r=s+1}^d \frac{1}{n_r} M_{n_r}^{[p_r]} \right), & \text{if } s > t, \end{cases}$$

where the matrices  $M_n^{[p]}$ ,  $H_n^{[p]}$ ,  $K_n^{[p]}$  are defined for all  $p, n \geq 1$  as follows:

$$\begin{aligned}
nK_n^{[p]} &:= \left[ \int_0^1 N'_{j+1,[p]} N'_{i+1,[p]} \right]_{i,j=1}^{n+p-2}, \\
H_n^{[p]} &:= \left[ \int_0^1 N'_{j+1,[p]} N_{i+1,[p]} \right]_{i,j=1}^{n+p-2}, \\
\frac{1}{n} M_n^{[p]} &:= \left[ \int_0^1 N_{j+1,[p]} N_{i+1,[p]} \right]_{i,j=1}^{n+p-2}.
\end{aligned}$$

From the analysis in [11], see in particular [11, Theorem 7, Eqs. (48)–(50) and Eqs. (65)–(68)], we know that  $M_n^{[p]}$ ,  $H_n^{[p]}$ ,  $K_n^{[p]}$  are small-rank perturbations of the Toeplitz matrices  $T_{n+p-2}(h_p)$ ,  $i T_{n+p-2}(g_p)$ ,  $T_{n+p-2}(f_p)$ , respectively, where the functions  $h_p$ ,  $g_p$ ,  $f_p$  are defined in (4.2)–(4.4). More precisely,

$$\begin{aligned}\text{rank}(M_n^{[p]} - T_{n+p-2}(h_p)) &\leq c_p, \\ \text{rank}(H_n^{[p]} - i T_{n+p-2}(g_p)) &\leq c_p, \\ \text{rank}(K_n^{[p]} - T_{n+p-2}(f_p)) &\leq c_p,\end{aligned}$$

with  $c_p$  a constant depending only on  $p$ ; for instance, we can take  $c_p = 4p - 2$  or  $c_p = 4p - 4$ , according to either [11, Eqs. (66), (68)] or [10, Eqs. (4.84), (4.86)]. Now we observe that, in the case  $s = t$ ,

$$\frac{\nu_s^2}{N(\mathbf{v})} T_{n+p-2}((H_p)_{ss}) = \frac{\nu_s^2}{N(\mathbf{v})} T_{n+p-2}(h_{p_1} \otimes \cdots \otimes h_{p_{s-1}} \otimes f_{p_s} \otimes h_{p_{s+1}} \otimes \cdots \otimes h_{p_d}),$$

and, recalling the relation  $\mathbf{n} = \mathbf{v}\mathbf{n}$ , we have

$$n^{d-2} \mathcal{L}_n^{[p]}(E_{ss}) = \frac{\nu_s^2}{N(\mathbf{v})} M_{n_1}^{[p_1]} \otimes \cdots \otimes M_{n_{s-1}}^{[p_{s-1}]} \otimes K_{n_s}^{[p_s]} \otimes M_{n_{s+1}}^{[p_{s+1}]} \otimes \cdots \otimes M_{n_d}^{[p_d]}.$$

Hence, by invoking Lemma 2.2 and the property (2.4), we see that the rank of the difference (4.10) in the case  $s = t$  is bounded from above by

$$N(\mathbf{n} + \mathbf{p} - \mathbf{2}) \sum_{i=1}^d \frac{c_{p_i}}{n_i + p_i - 2} = O(n^{d-1}).$$

The same argument shows that the rank of the difference (4.10) is  $O(n^{d-1})$  even in the case  $s \neq t$ .

Let us now prove (4.11). This is actually a direct consequence of the result that we have just proved and of the definition of sLT sequences. Nevertheless, since the latter definition is rather difficult, we include the details of the proof in order to make the reader more familiar with manipulations of sLT (and GLT) sequences. Because  $(H_p)_{st}$  is a separable trigonometric polynomial (see Section 2.4), the Toeplitz sequence  $\{(\nu_s \nu_t)/N(\mathbf{v}) T_{n+p-2}((H_p)_{st})\}_n$  is an sLT sequence. More precisely,

$$\left\{ \frac{\nu_s \nu_t}{N(\mathbf{v})} T_{n+p-2}((H_p)_{st}) \right\}_n \sim_{\text{sLT}} \left( \mathbf{1}, \frac{\nu_s \nu_t}{N(\mathbf{v})} (H_p)_{st} \right);$$

see [20, Theorem 5.1]. Therefore, by the definition of sLT sequences, for all sufficiently large  $\mathbf{m} \in \mathbb{N}^d$  there exists  $n_{\mathbf{m}}$  such that, for every  $n \geq n_{\mathbf{m}}$ ,

$$\frac{\nu_s \nu_t}{N(\mathbf{v})} T_{n+p-2}((H_p)_{st}) = \text{LT}_{n+p-2}^{\mathbf{m}} \left( \mathbf{1}, \frac{\nu_s \nu_t}{N(\mathbf{v})} (H_p)_{st} \right) + R_{n+p-2, \mathbf{m}} + S_{n+p-2, \mathbf{m}}, \quad (4.12)$$

where the first matrix in the right-hand side is defined in [20, Definition 2.1],

$$\text{rank}(R_{n+p-2, \mathbf{m}}) \leq c(\mathbf{m}) N(\mathbf{n} + \mathbf{p} - \mathbf{2}) \sum_{j=1}^d (n_j + p_j - 2)^{-1},$$

$$\|S_{n+p-2, \mathbf{m}}\|_1 \leq \omega(\mathbf{m}) N(\mathbf{n} + \mathbf{p} - \mathbf{2}).$$

In the above inequalities, the quantities  $c(\mathbf{m})$ ,  $\omega(\mathbf{m})$  depend on  $\mathbf{m}$  but not on  $n$ , and  $\lim_{\mathbf{m} \rightarrow \infty} \omega(\mathbf{m}) = 0$ . The rank of the difference (4.10) is  $O(n^{d-1})$ , so it can be bounded by a constant  $C$  (independent of  $n$ ) times  $N(\mathbf{n} + \mathbf{p} - \mathbf{2}) \sum_{j=1}^d (n_j + p_j - 2)^{-1}$ . Therefore, the matrix  $n^{d-2} \mathcal{L}_n^{[p]}(E_{st})$  admits the same decomposition as (4.12), with the only difference that  $R_{n+p-2, \mathbf{m}}$  is replaced by another small-rank term  $\widetilde{R}_{n+p-2, \mathbf{m}}$  and  $c(\mathbf{m})$  is replaced by  $c(\mathbf{m}) + C$ . It follows that  $\{n^{d-2} \mathcal{L}_n^{[p]}(E_{st})\}_n$  is an sLT sequence and that (4.11) holds; see [20, Definition 2.1]. This concludes the proof of Step 3.

## 4.5 Proof of Step 4

Let  $1 \leq s, t \leq d$ . We show that, whenever  $a : [0, 1]^d \rightarrow \mathbb{R}$  is continuous,

$$\|n^{d-2} \mathcal{L}_n^{[p]}(a(\hat{\mathbf{x}})E_{st}) - n^{d-2} D_{n+p-2}(a) \mathcal{L}_n^{[p]}(E_{st})\| = O(\omega_a(n^{-1})), \quad (4.13)$$

and

$$\{n^{d-2} \mathcal{L}_n^{[p]}(a(\hat{\mathbf{x}})E_{st})\}_n \sim_{\text{sLT}} \left( a, \frac{\nu_s \nu_t}{N(\mathbf{v})} (H_p)_{st} \right). \quad (4.14)$$

In the following,  $C$  denotes a generic constant independent of  $n$ .

For every  $\mathbf{i}, \mathbf{j} = \mathbf{2}, \dots, \mathbf{n} + \mathbf{p} - \mathbf{1}$ ,

$$\begin{aligned} \left| (\mathcal{L}_n^{[p]}(a(\hat{\mathbf{x}})E_{st}) - D_{n+p-2}(a) \mathcal{L}_n^{[p]}(E_{st}))_{i-1, j-1} \right| &= \left| (\mathcal{L}_n^{[p]}(a(\hat{\mathbf{x}})E_{st}))_{i-1, j-1} - (D_{n+p-2}(a))_{i-1, i-1} (\mathcal{L}_n^{[p]}(E_{st}))_{i-1, j-1} \right| \\ &= \left| \int_{[0,1]^d} \left[ a(\hat{\mathbf{x}}) - a\left(\frac{\mathbf{i}-\mathbf{1}}{\mathbf{n}+\mathbf{p}-\mathbf{2}}\right) \right] \frac{\partial N_{j, [p]}}{\partial \hat{x}_s} \frac{\partial N_{i, [p]}}{\partial \hat{x}_t} \right| \\ &\leq \max_{\hat{\mathbf{x}} \in \text{supp}(N_{i, [p]})} \left| a(\hat{\mathbf{x}}) - a\left(\frac{\mathbf{i}-\mathbf{1}}{\mathbf{n}+\mathbf{p}-\mathbf{2}}\right) \right| \int_{[0,1]^d} \left| \frac{\partial N_{j, [p]}}{\partial \hat{x}_s} \right| \left| \frac{\partial N_{i, [p]}}{\partial \hat{x}_t} \right|. \end{aligned} \quad (4.15)$$

Using the relations  $\text{supp}(N_{i, [p]}) = [t_i, t_{i+p+1}] \times \dots \times [t_d, t_{d+p_d+1}]$  and

$$\left| \hat{x} - \frac{i-1}{n+p-2} \right| \leq Cn^{-1}, \quad \forall \hat{x} \in \text{supp}(N_{i, [p]}) = [t_i, t_{i+p+1}],$$

we see that

$$\max_{\hat{\mathbf{x}} \in \text{supp}(N_{i, [p]})} \left\| \hat{\mathbf{x}} - \frac{\mathbf{i}-\mathbf{1}}{\mathbf{n}+\mathbf{p}-\mathbf{2}} \right\|_{\infty} \leq C \left( \min_{r=1, \dots, d} n_r \right)^{-1} \leq Cn^{-1}.$$

Therefore, taking into account the definition of the modulus of continuity

$$\omega_a(\delta) := \max_{\substack{\hat{\mathbf{x}}, \hat{\mathbf{y}} \in [0,1]^d \\ \|\hat{\mathbf{x}} - \hat{\mathbf{y}}\|_{\infty} \leq \delta}} |a(\hat{\mathbf{x}}) - a(\hat{\mathbf{y}})|, \quad \delta \geq 0,$$

we obtain

$$\max_{\hat{\mathbf{x}} \in \text{supp}(N_{i, [p]})} \left| a(\hat{\mathbf{x}}) - a\left(\frac{\mathbf{i}-\mathbf{1}}{\mathbf{n}+\mathbf{p}-\mathbf{2}}\right) \right| \leq C\omega_a(n^{-1}).$$

Moreover, we know from the proof of Step 1 that

$$\int_{[0,1]^d} \left| \frac{\partial N_{j, [p]}}{\partial \hat{x}_s} \right| \left| \frac{\partial N_{i, [p]}}{\partial \hat{x}_t} \right| \leq Cn^{-d+2}.$$

Hence, from (4.15) we obtain

$$\left| (\mathcal{L}_n^{[p]}(a(\hat{\mathbf{x}})E_{st}) - D_{n+p-2}(a) \mathcal{L}_n^{[p]}(E_{st}))_{i-1, j-1} \right| \leq C\omega_a(n^{-1})n^{-d+2}$$

and

$$\left| (n^{d-2} \mathcal{L}_n^{[p]}(a(\hat{\mathbf{x}})E_{st}) - n^{d-2} D_{n+p-2}(a) \mathcal{L}_n^{[p]}(E_{st}))_{i-1, j-1} \right| \leq C\omega_a(n^{-1}),$$

for every  $\mathbf{i}, \mathbf{j} = \mathbf{2}, \dots, \mathbf{n} + \mathbf{p} - \mathbf{1}$ . In addition, the matrix  $n^{d-2} \mathcal{L}_n^{[p]}(a(\hat{\mathbf{x}})E_{st}) - n^{d-2} D_{n+p-2}(a) \mathcal{L}_n^{[p]}(E_{st})$ , like the matrices  $R_{\mathbf{G}, n}^{[p]}$  and  $K_{\mathbf{G}, n}^{[p]}$  already considered in the proof of Step 1, has the entry 0 in each position  $(\mathbf{i}, \mathbf{j})$  whenever  $\|\mathbf{i} - \mathbf{j}\|_{\infty} \geq \|\mathbf{p}\|_{\infty} + 1$ , due to the local support of the B-splines. It follows that the number of non-zero components of  $n^{d-2} \mathcal{L}_n^{[p]}(a(\hat{\mathbf{x}})E_{st}) - n^{d-2} D_{n+p-2}(a) \mathcal{L}_n^{[p]}(E_{st})$  in each row and column is at most  $(2\|\mathbf{p}\|_{\infty} + 1)^d$ , and so

$$\left\| n^{d-2} \mathcal{L}_n^{[p]}(a(\hat{\mathbf{x}})E_{st}) - n^{d-2} D_{n+p-2}(a) \mathcal{L}_n^{[p]}(E_{st}) \right\| \leq C\omega_a(n^{-1}).$$

This proves (4.13).

For the proof of (4.14), let us denote by  $O$  the zero matrix, whose size will be clear from the context. By examining the proof of [20, Theorem 5.1], we see that the splitting (4.12) holds with  $S_{n+p-2,m} = O$ . When proving Step 3, we showed that  $n^{d-2} \mathcal{L}_n^{[p]}(E_{st})$  admits the same decomposition as (4.12) with  $R_{n+p-2,m}$  replaced by  $\tilde{R}_{n+p-2,m}$  and with  $c(\mathbf{m})$  replaced by  $c(\mathbf{m}) + C$ , so the small-norm term  $S_{n+p-2,m}$  equals  $O$  also for  $n^{d-2} \mathcal{L}_n^{[p]}(E_{st})$ . Concerning the diagonal sampling matrix  $D_{n+p-2}(a)$ , it follows from the definition of sLT sequences that  $\{D_{n+p-2}(a)\}_n \sim_{\text{sLT}} (a, 1)$ . More precisely, for all sufficiently large  $\mathbf{m} \in \mathbb{N}^d$  there exists  $n_m$  such that, for every  $n \geq n_m$ ,

$$D_{n+p-2}(a) = \text{LT}_{n+p-2}^m(a, 1) + \tilde{R}_{n+p-2,m} + \tilde{S}_{n+p-2,m},$$

where the first matrix in the right-hand side is defined in [20, Definition 2.1] and

$$\begin{aligned} \text{rank}(\tilde{R}_{n+p-2,m}) &\leq \check{c}(\mathbf{m})N(\mathbf{n} + \mathbf{p} - \mathbf{2}) \sum_{j=1}^d (n_j + p_j - 2)^{-1}, \\ \|\tilde{S}_{n+p-2,m}\| &\leq \omega_a \left( (\min_{j=1,\dots,d} m_j)^{-1} \right). \end{aligned}$$

We also refer the reader to [20, Example 5.4].<sup>3</sup> Then, by [20, Theorem 5.3], we have

$$\{n^{d-2} D_{n+p-2}(a) \mathcal{L}_n^{[p]}(E_{st})\}_n \sim_{\text{sLT}} \left( a, \frac{v_s v_t}{N(\mathbf{v})} (H_p)_{st} \right).$$

Using this relation together with (4.13), and applying again the definition of sLT sequences, we get (4.14).

Finally, we note that in Steps 5–6 there is nothing to prove.

#### 4.6 Final remarks

Assume that  $K := [\kappa_{ij}]_{i,j=1}^d$  and  $u$  are sufficiently regular, say  $\kappa_{ij} \in C^1(\Omega) \cap C(\overline{\Omega})$  for all  $i, j = 1, \dots, d$ , and  $u \in C^2(\Omega) \cap C(\overline{\Omega})$ . In this case, our problem (1.1) can be reformulated as follows:

$$\begin{cases} -1(K \circ Hu) \mathbf{1}^T + \boldsymbol{\beta} \cdot \nabla u + \gamma u = f, & \text{in } \Omega, \\ u = 0, & \text{on } \partial, \end{cases} \quad (4.16)$$

where  $Hu$  is the Hessian of  $u$ ,

$$(Hu)_{ij} := \frac{\partial^2 u}{\partial x_i \partial x_j},$$

and  $\boldsymbol{\beta}$  collects the coefficients of the first-order derivatives in (1.1):

$$\beta_j := \alpha_j - \sum_{i=1}^d \frac{\partial \kappa_{ij}}{\partial x_i}.$$

For any  $u : \overline{\Omega} \rightarrow \mathbb{R}$ , consider the corresponding function defined on the parametric domain  $[0, 1]^d$  by

$$\hat{u} : [0, 1]^d \rightarrow \mathbb{R}, \quad \hat{u}(\hat{\mathbf{x}}) := u(\mathbf{x}), \quad \mathbf{x} = \mathbf{G}(\hat{\mathbf{x}}).$$

In other words,  $\hat{u} := u(\mathbf{G})$ . Then,  $u$  satisfies (4.16) if and only if  $\hat{u}$  satisfies the corresponding transformed problem

$$\begin{cases} -1(K_G \circ H\hat{u}) \mathbf{1}^T + \boldsymbol{\beta}_G \cdot \nabla \hat{u} + \gamma(\mathbf{G}) \hat{u} = f(\mathbf{G}), & \text{in } (0, 1)^d, \\ \hat{u} = 0, & \text{on } \partial((0, 1)^d), \end{cases} \quad (4.17)$$

<sup>3</sup>Take into account that in [20, Example 5.4] the index  $n = (n_1, \dots, n_d)$  is a  $d$ -index and the matrix  $D_{n,a}$  coincides with our diagonal sampling matrix  $D_n(a)$ , as defined in (2.5) for  $\mathbf{m} = n$ .

where  $H\hat{u}$  is the Hessian of  $\hat{u}$ ,

$$(H\hat{u})_{ij} = \frac{\partial^2 \hat{u}}{\partial \hat{x}_i \partial \hat{x}_j},$$

$K_G$  is given in (3.5), and  $\beta_G$  is the transformed advection coefficient of the PDE, whose expression in terms of  $K$ ,  $\beta$ ,  $\mathbf{G}$  is complicated and hence not reported here.

Now assume that  $\Omega = (0, 1)^d$  and take  $\mathbf{G}$  equal to the identity map. In this case, problems (4.16)–(4.17) are the same,  $\mathbf{x} = \hat{\mathbf{x}}$ ,  $u = \hat{u}$ ,  $K_G = K$ , and the symbol (4.5) reduces to

$$\frac{\boldsymbol{\nu}(K(\hat{\mathbf{x}}) \circ H_p(\boldsymbol{\theta}))\boldsymbol{\nu}^T}{N(\boldsymbol{\nu})}.$$

In particular, if we give to every direction the same attention by choosing  $\boldsymbol{\nu} = \mathbf{1}$ , then the symbol is

$$\mathbf{1}(K(\hat{\mathbf{x}}) \circ H_p(\boldsymbol{\theta}))\mathbf{1}^T.$$

It is quite remarkable that the main operator in the Hörmander sense [16], namely  $-\mathbf{1}(K \circ H\hat{u})\mathbf{1}^T$ , has a discrete spectral counterpart  $\mathbf{1}(K(\hat{\mathbf{x}}) \circ H_p(\boldsymbol{\theta}))\mathbf{1}^T$  which looks formally the same. The similarity becomes even more evident if we recall (from Step 3 in Section 4.2) that  $H_p(\boldsymbol{\theta})$  is a matrix of trigonometric polynomials in the Fourier variables, whose components are related to the discretization of the components of

$$-H := \left[ -\frac{\partial^2}{\partial \hat{x}_i \partial \hat{x}_j} \right]_{i,j=1}^d,$$

i.e., the opposite of the Hessian operator in the parametric variables  $\hat{\mathbf{x}}$ .<sup>4</sup>

In the case of an arbitrary domain  $\Omega$  described by a non-trivial geometry map  $\mathbf{G}$ , the symbol is

$$\frac{\boldsymbol{\nu}(|\det(J_G(\hat{\mathbf{x}}))| K_G(\hat{\mathbf{x}}) \circ H_p(\boldsymbol{\theta}))\boldsymbol{\nu}^T}{N(\boldsymbol{\nu})},$$

and when  $\boldsymbol{\nu} = \mathbf{1}$  it reduces to

$$\mathbf{1}(|\det(J_G(\hat{\mathbf{x}}))| K_G(\hat{\mathbf{x}}) \circ H_p(\boldsymbol{\theta}))\mathbf{1}^T.$$

Even in this case, the symbol preserves the formal structure of the main differential operator  $-\mathbf{1}(K_G \circ H\hat{u})\mathbf{1}^T$  associated with the problem (4.17). We see, however, the appearance of the determinant factor  $|\det(J_G)|$ . This factor is not present when our PDE is approximated by isogeometric collocation methods [8], in which case the resulting symbol is given by

$$\boldsymbol{\nu}(K_G(\hat{\mathbf{x}}) \circ H_p(\boldsymbol{\theta}))\boldsymbol{\nu}^T,$$

and is formally identical to the main operator  $-\mathbf{1}(K_G \circ H\hat{u})\mathbf{1}^T$ , especially when  $\boldsymbol{\nu} = \mathbf{1}$ . The determinant factor  $|\det(J_G)|$  appearing in the context of isogeometric Galerkin methods is interpreted by saying that the symbol of the (normalized) isogeometric Galerkin matrices is less ill-conditioned than the symbol of the (normalized) isogeometric collocation matrices. Indeed, since  $|\det(J_G)|$  appears in the numerator of the symbol, it helps in keeping the conditioning of the IgA Galerkin matrices moderate when the geometry map  $\mathbf{G}$  is (nearly) singular. A geometry map  $\mathbf{G}$  is (nearly) singular when  $\det(J_G)$  is (nearly) zero at one or more points. We refer the reader to [22, Section 5.2] for the analysis of the map effect on the conditioning in the collocation setting.

---

<sup>4</sup>The matrix  $H_p(\boldsymbol{\theta})$  is sometimes referred to as ‘the symbol of the (negative) Hessian operator’. With a careful consideration of the results presented in this paper, the reader could guess the origin of this terminology, which, however, is not completely rigorous from the mathematical viewpoint.

We end this section with the following observation. Assume that  $\Omega = (0, 1)^d$ ,  $\mathbf{G}$  is the identity map, and  $K = I$  is the identity matrix. In this case, the function  $f_{\mathbf{G}, p}^{(v)}$  is independent of the  $\hat{\mathbf{x}}$ -variables. Therefore, we can consider its restriction to the domain  $[-\pi, \pi]^d$ , i.e., the function  $f_p^{(v)} : [-\pi, \pi]^d \rightarrow \mathbb{R}$ ,

$$f_p^{(v)}(\boldsymbol{\theta}) := \frac{1}{N(\mathbf{v})} \sum_{i=1}^d v_i^2 (h_{p_i} \otimes \cdots \otimes h_{p_{i-1}} \otimes f_{p_i} \otimes h_{p_{i+1}} \otimes \cdots \otimes h_{p_d})(\boldsymbol{\theta}). \quad (4.18)$$

A direct application of Definition 2.1 shows that  $f_p^{(v)}$  is an eigenvalue and singular value symbol for the normalized sequence  $\{n^{d-2} A_n^{[p]}\}_n$ , where  $\mathbf{n} = \mathbf{v}n$  and  $A_n^{[p]}$  is the matrix resulting from the Galerkin B-spline IgA approximation of (1.1) when  $\Omega = (0, 1)^d$ ,  $\mathbf{G}$  is the identity map, and  $K = I$ . The symbol  $f_p^{(v)}$  is exactly the one appearing in [10, Chapter 4], where problem (1.1) was considered over the hypercube  $(0, 1)^d$  in the constant coefficient case ( $K = I$ ); see also [11] for the original analysis in the univariate and bivariate settings.

## 5 Properties of the symbol

The formal structure of the symbol  $f_{\mathbf{G}, p}^{(v)}$ , discussed in the previous section, is actually not surprising, because it is analogous to the structure of the symbols arising in the context of other approximation techniques, such as FDs and FEs. The determinant factor  $|\det(J_{\mathbf{G}})|$ , appearing in the context of isogeometric Galerkin methods, was already observed in the FE setting (see [1]), but it is not present in the framework of isogeometric collocation methods (see [8]) and FDs (see [20, Section 6] and [22]). Although the formal structure of the symbol  $f_{\mathbf{G}, p}^{(v)}$  is shared by different approximation techniques, certain analytic features are not so common. In this section, we investigate the properties of the symbol  $f_{\mathbf{G}, p}^{(v)}$ , for which we heavily rely on results obtained in [8].

We first recall that the functions  $h_p, g_p, f_p$  in (4.2)–(4.4) coincide with the functions  $h_{2p+1}, g_{2p+1}, f_{2p+1}$  defined in [8, Eqs. (3.7)–(3.9)], which have been deeply studied in [8]. From the properties of these functions, we get the following results. The first one concerns the symbol  $f_p^{(v)}$  introduced in (4.18), and can be deduced from [8, Lemmas 3.4 and 3.6].

**Theorem 5.1.** *The following properties hold.*

1.  $c_{p, v} \sum_{k=1}^d (2 - 2 \cos \theta_k) \leq f_p^{(v)}(\boldsymbol{\theta}) \leq C_v \sum_{k=1}^d (2 - 2 \cos \theta_k)$ , where

$$c_{p, v} := \left( \frac{4}{\pi^2} \right)^{\sum_{i=1}^d p_i + d - 1} \frac{\min_{i=1, \dots, d} v_i^2}{N(\mathbf{v})}, \quad C_v := \frac{\max_{i=1, \dots, d} v_i^2}{N(\mathbf{v})}.$$

2. Let  $M_{f_p^{(v)}} := \max_{\boldsymbol{\theta} \in [-\pi, \pi]^d} f_p^{(v)}(\boldsymbol{\theta})$ , then

$$\frac{f_p^{(v)}(\theta_1, \dots, \theta_{j-1}, \pi, \theta_{j+1}, \dots, \theta_d)}{M_{f_p^{(v)}}} \leq \frac{f_p^{(v)}(\theta_1, \dots, \theta_{j-1}, \pi, \theta_{j+1}, \dots, \theta_d)}{f_p^{(v)}(\theta_1, \dots, \theta_{j-1}, \frac{\pi}{2}, \theta_{j+1}, \dots, \theta_d)} \leq 2^{2-p_j},$$

for all  $j = 1, \dots, d$ .

In particular,  $f_p^{(v)}$  has a unique zero of order two at  $\boldsymbol{\theta} = \mathbf{0}$ , like the function  $\sum_{k=1}^d (2 - 2 \cos \theta_k)$ , and the value  $f_p^{(v)}(\theta_1, \dots, \theta_{j-1}, \pi, \theta_{j+1}, \dots, \theta_d) / M_{f_p^{(v)}}$  converges to 0 exponentially when  $p_j \rightarrow \infty$ .

The zero of the symbol  $f_p^{(v)}$  at  $\theta = \mathbf{0}$  is expected, because it is a canonical feature of the symbol associated with the discretization matrices of differential problems like (1.1). Indeed, it was already observed in the FD and FE settings. Such a zero is interpreted as a source of ill-conditioning for the corresponding IgA Galerkin matrices  $n^{d-2}A_n^{[p]}$  in the low frequencies; see [6, Section 3.2.2] for the terminology of frequencies. On the other hand, when the spline degrees  $p$  are large, the normalized symbol  $f_p^{(v)}/M_{f_p^{(v)}}$  takes very small values at the so-called  $\pi$ -edge points

$$\{\theta \in [0, \pi]^d : \theta_i = \pi \text{ for some } i\}. \quad (5.1)$$

These ‘numerical zeros’ are interpreted as an ill-conditioning of  $n^{d-2}A_n^{[p]}$  in the high frequencies corresponding to such points. This second (non-canonical) source of ill-conditioning is responsible for the convergence deterioration of standard multigrid methods when the  $p_i$  increase. A way to overcome this problem consists in adopting a multi-iterative strategy, as suggested in [6] (see also [7]).

The next result concerns the matrix  $H_p$ . We use the abbreviation SPSD for Symmetric Positive Semi-Definite.

**Theorem 5.2.** *For  $d = 1, 2, 3$ , the matrix  $H_p(\theta)$  in (4.1) is SPSD for all  $\theta \in [-\pi, \pi]^d$  and SPD for all  $\theta \in [-\pi, \pi]^d$  such that  $\theta_1 \cdots \theta_d \neq 0$ .*

*Proof.* The result for  $d = 1, 2$  follows immediately from [8, Lemma 3.6 and Theorem 5.2], taking into account that  $H_p = f_p$  for  $d = 1$ . In the remainder, we focus on the case  $d = 3$ .

The symmetry of  $H_p = H_{p_1, p_2, p_3}$  can be directly seen from its definition, i.e.,

$$H_{p_1, p_2, p_3} := \begin{bmatrix} f_{p_1} \otimes h_{p_2} \otimes h_{p_3} & g_{p_1} \otimes g_{p_2} \otimes h_{p_3} & g_{p_1} \otimes h_{p_2} \otimes g_{p_3} \\ g_{p_1} \otimes g_{p_2} \otimes h_{p_3} & h_{p_1} \otimes f_{p_2} \otimes h_{p_3} & h_{p_1} \otimes g_{p_2} \otimes g_{p_3} \\ g_{p_1} \otimes h_{p_2} \otimes g_{p_3} & h_{p_1} \otimes g_{p_2} \otimes g_{p_3} & h_{p_1} \otimes h_{p_2} \otimes f_{p_3} \end{bmatrix}.$$

For every  $p \geq 1$ , let  $e_p := f_p h_p - (g_p)$ . Then,

$$\begin{aligned} \det(H_{p_1, p_2, p_3}) &= (f_{p_1} \otimes h_{p_2} \otimes h_{p_3})(h_{p_1} \otimes f_{p_2} \otimes h_{p_3})(h_{p_1} \otimes h_{p_2} \otimes f_{p_3}) \\ &\quad + 2(h_{p_1} \otimes g_{p_2} \otimes g_{p_3})(g_{p_1} \otimes h_{p_2} \otimes g_{p_3})(g_{p_1} \otimes g_{p_2} \otimes h_{p_3}) \\ &\quad - (f_{p_1} \otimes h_{p_2} \otimes h_{p_3})(h_{p_1} \otimes g_{p_2} \otimes g_{p_3})(h_{p_1} \otimes g_{p_2} \otimes g_{p_3}) \\ &\quad - (g_{p_1} \otimes h_{p_2} \otimes g_{p_3})(h_{p_1} \otimes f_{p_2} \otimes h_{p_3})(g_{p_1} \otimes h_{p_2} \otimes g_{p_3}) \\ &\quad - (g_{p_1} \otimes g_{p_2} \otimes h_{p_3})(g_{p_1} \otimes g_{p_2} \otimes h_{p_3})(h_{p_1} \otimes h_{p_2} \otimes f_{p_3}) \\ &= (h_{p_1} \otimes h_{p_2} \otimes h_{p_3}) [(e_{p_1} + (g_{p_1})^2) \otimes (e_{p_2} + (g_{p_2})^2) \otimes (e_{p_3} + (g_{p_3})^2) \\ &\quad - e_{p_1} \otimes (g_{p_2})^2 \otimes (g_{p_3})^2 - (g_{p_1})^2 \otimes e_{p_2} \otimes (g_{p_3})^2 \\ &\quad - (g_{p_1})^2 \otimes (g_{p_2})^2 \otimes (e_{p_3} + (g_{p_3})^2)] \\ &= (h_{p_1} \otimes h_{p_2} \otimes h_{p_3}) [e_{p_1} \otimes e_{p_2} \otimes e_{p_3} + (g_{p_1})^2 \otimes e_{p_2} \otimes e_{p_3} \\ &\quad + e_{p_1} \otimes (g_{p_2})^2 \otimes e_{p_3} + e_{p_1} \otimes e_{p_2} \otimes (g_{p_3})^2]. \end{aligned}$$

We recall from [8, Lemmas 3.4–3.7] that

$$\begin{aligned} h_p(\theta) &> 0, & \theta &\in [-\pi, \pi], \\ f_p(\theta) &> 0, & \theta &\in [-\pi, \pi] \setminus \{0\}, \\ e_p(\theta) &> 0, & \theta &\in [-\pi, \pi] \setminus \{0\}, \end{aligned}$$

and

$$f_p(0) = 0, \quad g_p(0) = 0, \quad e_p(0) = 0.$$

Hence, we deduce that  $\det(H_{p_1, p_2, p_3}) \geq 0$  over  $[-\pi, \pi]^3$ . Moreover, if  $\theta_1 \theta_2 \theta_3 \neq 0$  then  $\det(H_{p_1, p_2, p_3}) > 0$ . In addition, from [8, Lemma 3.6 and Theorem 5.2] we infer that  $(H_{p_1, \dots, p_d})_{1,1} > 0$  if  $\theta_1 \neq 0$  and  $\det([(H_{p_1, p_2, p_3})_{i,j}]_{i,j=1}^2) > 0$  if  $\theta_1 \theta_2 \neq 0$ . Therefore, Sylvester's criterion implies that  $H_{p_1, p_2, p_3}$  is SPD if  $\theta_1 \theta_2 \theta_3 \neq 0$ . Moreover,  $H_{p_1, p_2, p_3}$  is SPSD over  $[-\pi, \pi]^3$  by a continuity argument.  $\square$

*Remark 5.3.* We conjecture that Theorem 5.2 holds for every  $d \geq 1$ . It is clear from its definition (4.1) that  $H_p$  is symmetric for any  $d$ . To prove the positive definiteness of  $H_p$ , it suffices to prove that the determinant of  $H_p$  is positive for any  $d$ . Indeed, from (4.1) we see that the upper-left  $k \times k$  submatrix of  $H_p = H_{p_1, \dots, p_d}$  can be written as

$$[(H_{p_1, \dots, p_d})_{i,j}]_{i,j=1}^k = [(H_{p_1, \dots, p_k})_{i,j} \otimes (\bigotimes_{r=k+1}^d h_{p_r})]_{i,j=1}^k,$$

and so

$$\det([(H_{p_1, \dots, p_d})_{i,j}]_{i,j=1}^k) = \det(H_{p_1, \dots, p_k}) \otimes (\bigotimes_{r=k+1}^d (h_{p_r})^k).$$

Since  $h_p > 0$  over  $[-\pi, \pi]$ , the  $k$ -th leading principal minor of  $H_p$  is positive if  $\det(H_{p_1, \dots, p_k}) > 0$ , and  $H_p$  is positive definite if

$$\det(H_{p_1, \dots, p_k}) > 0, \quad k = 1, \dots, d.$$

The fact that  $H_p$  is SPD follows by induction from Sylvester's criterion.

In the following, we assume that  $H_p(\theta)$  is SPSD for all  $\theta \in [-\pi, \pi]^d$ . Under this assumption, we derive some interesting properties of the symbol  $f_{G,p}^{(v)}$ , which are certainly true for  $d = 1, 2, 3$ , by Theorem 5.2.

**Theorem 5.4.** *Assume that  $H_p(\theta)$  is SPSD for all  $\theta \in [-\pi, \pi]^d$ . Then, the following properties hold.*

1. The symbol  $f_{G,p}^{(v)}$  is non-negative a.e. in  $[0, 1]^d \times [-\pi, \pi]^d$ .
2. For every  $(\hat{x}, \theta) \in [0, 1]^d \times [-\pi, \pi]^d$ , we have

$$\begin{aligned} f_{G,p}^{(v)}(\hat{x}, \theta) &\geq \lambda_{\min}(K_G(\hat{x})) |\det(J_G(\hat{x}))| f_p^{(v)}(\theta), \\ f_{G,p}^{(v)}(\hat{x}, \theta) &\leq \lambda_{\max}(K_G(\hat{x})) |\det(J_G(\hat{x}))| f_p^{(v)}(\theta). \end{aligned} \tag{5.2}$$

In particular, if

$$cI \leq K_G |\det(J_G)| \leq CI \quad \text{a.e. in } [0, 1]^d, \quad \text{for some } c, C > 0, \tag{5.3}$$

then

$$c f_p^{(v)}(\theta) \leq f_{G,p}^{(v)}(\hat{x}, \theta) \leq C f_p^{(v)}(\theta),$$

for all  $\theta \in [-\pi, \pi]^d$  and for almost every  $\hat{x} \in [0, 1]^d$ .

*Proof.* The first statement can be shown by following the argument used in the proof of [8, Theorem 5.4] or [10, Theorem 5.5], also taking into account that  $K$  is SPD over  $\Omega$  by hypothesis.

Let us now prove the second statement. From the properties of the Hadamard product (see e.g. the proof of [8, Theorem 5.2] or [10, Lemma 1.6]), we know that  $X \circ Y \geq Z \circ Y$  whenever  $Y$  is Hermitian positive semi-definite and  $X \geq Z$ . Moreover, it is clear that

$$\lambda_{\min}(K_G(\hat{x}))I \leq K_G(\hat{x}) \leq \lambda_{\max}(K_G(\hat{x}))I.$$

Therefore, by using the assumption that  $H_p(\theta)$  is SPSD and the definitions of  $f_{G,p}^{(v)}$  and  $f_p^{(v)}$ , we get (5.2).  $\square$

*Remark 5.5.* Alternatively, we could prove the first statement in Theorem 5.4 even without the assumption that  $H_p(\boldsymbol{\theta})$  is SPSD for all  $\boldsymbol{\theta} \in [-\pi, \pi]^d$ , at least in the case where  $\mathbf{G}$  is regular and the components of  $K$  are continuous over  $\overline{\Omega}$ . Indeed, we know from Theorem 4.1 (applied with  $\boldsymbol{\alpha} = \mathbf{0}$  and  $\gamma = 0$ ) that  $\{n^{d-2}K_{\mathbf{G},n}^{[p]}\}_n \sim_\lambda f_{\mathbf{G},p}^{(v)}$ . Therefore, by [15, Theorem 2.4], each point of the essential range of  $f_{\mathbf{G},p}^{(v)}$  strongly attracts the spectrum of  $n^{d-2}K_{\mathbf{G},n}^{[p]}$  with infinite order (see [15, Definition 2.3] for the concept of spectral attraction). Roughly speaking, this means that, for each point in the essential range of  $f_{\mathbf{G},p}^{(v)}$ , the number of eigenvalues of  $n^{d-2}K_{\mathbf{G},n}^{[p]}$  that collapse on this point tends to  $\infty$  with  $n$ . Since every matrix  $n^{d-2}K_{\mathbf{G},n}^{[p]}$  is positive definite, it is clear that no eigenvalue of  $n^{d-2}K_{\mathbf{G},n}^{[p]}$  can collapse on a negative real number. Hence, the essential range of  $f_{\mathbf{G},p}^{(v)}$  is contained in  $[0, \infty)$  and, consequently,  $f_{\mathbf{G},p}^{(v)} \geq 0$  (a.e.) over its domain  $[0, 1]^d \times [-\pi, \pi]^d$ .

We observe that the condition (5.3) is usually satisfied in practice. For instance, it is satisfied if

- $K \geq c_K I$  a.e. in  $\Omega$ , for some  $c_K > 0$ ;
- $\mathbf{G}$  is regular, like in Theorem 4.1.

Note that we only require the inequality  $K \geq c_K I$  and not the opposite  $K \leq C_K I$ , because the latter is certainly satisfied, since the components of  $K$  are assumed to be in  $L^\infty(\Omega)$ .

Assuming that (5.3) is met and  $H_p$  is SPSD over  $[-\pi, \pi]^d$ , the second statement of Theorem 5.4 implies that  $f_{\mathbf{G},p}^{(v)}$  has the same behavior of  $f_p^{(v)}$ . In particular, for (almost) every  $\hat{\mathbf{x}} \in [0, 1]^d$ , the function  $f_{\mathbf{G},p}^{(v)}(\hat{\mathbf{x}}, \cdot) : [-\pi, \pi]^d \rightarrow \mathbb{R}$  has a unique zero of order two at  $\boldsymbol{\theta} = \mathbf{0}$ , like the functions  $f_p^{(v)}$  and  $\sum_{k=1}^d (2 - 2 \cos \theta_k)$ . Moreover, if one of the  $p$  parameters, say  $p_i$ , is large, then  $f_{\mathbf{G},p}^{(v)}(\hat{\mathbf{x}}, \cdot)$  also has infinitely many numerical zeros located at the  $\pi$ -edge points (5.1).

## 6 Conclusions

We have presented an asymptotic spectral analysis of the discretization matrices associated with the Galerkin B-spline IgA approximation of full elliptic PDEs. In particular, we have computed the spectral and singular value symbol  $f_{\mathbf{G},p}^{(v)}$ , in the sense of Definition 2.1, of the normalized sequence of matrices  $\{n^{d-2}A_{\mathbf{G},n}^{[p]}\}_n$ , with  $\mathbf{n} = \nu n$ . We have also collected some properties of the symbol  $f_{\mathbf{G},p}^{(v)}$ , which – in agreement with previous results in the FD/FE/collocation contexts [1, 8, 20, 21] – has the canonical structure described in item b) of the introduction. Looking at the range of  $f_{\mathbf{G},p}^{(v)}$ , it has been observed that  $f_{\mathbf{G},p}^{(v)}$  is a non-negative function with a unique zero of order two at  $\boldsymbol{\theta} = \mathbf{0}$ . However, when  $p_i$  is large, the symbol also shows infinitely many ‘numerical zeros’ at the points  $(\hat{\mathbf{x}}, \boldsymbol{\theta})$  such that  $\theta_i = \pi$ . While the zero at  $\boldsymbol{\theta} = \mathbf{0}$  is expected, because it is common to any approximation method (see analogous features in the FD/FE/collocation symbols [1, 8, 20, 21]), the second type of zeros leads to the surprising fact that, for large  $\|\mathbf{p}\|_\infty$ , there is a subspace of high frequencies where the (normalized) Galerkin B-spline IgA stiffness matrices are ill-conditioned. This non-canonical feature is responsible for the slowdown, with respect to  $\mathbf{p}$ , of standard iterative methods. On the other hand, its knowledge and the knowledge of other properties of the (simplified) symbol  $f_p^{(v)}$  allowed us to construct a (multi-iterative) multigrid solver involving the PCG/P-GMRES as a smoother at the finest level, for which optimal convergence properties are numerically observed, with a remarkable robustness with respect to all the relevant parameters; see [6].

The pattern of the proof presented herein for computing the symbol  $f_{\mathbf{G},p}^{(v)}$  is based on the theory of sLT and GLT sequences [20, 21], as well as on the results in [12, 15] (for the non-symmetric case), and it is very general. A future line of research can include the application of this very general technique to the computation of the symbol of the discretization matrices associated with (1.1) by means of other numerical methods, such as, for example, the isogeometric Galerkin and collocation methods based on NURBS instead of B-splines.

## References

- [1] B. BECKERMANN, S. SERRA-CAPIZZANO. *On the asymptotic spectrum of finite element matrix sequences*. SIAM J. Numer. Anal. **45** (2007) 746–769.
- [2] R. BHATIA. *Matrix analysis*. Springer-Verlag, New York (1997).
- [3] C. DE BOOR. *A practical guide to splines*. Springer-Verlag, New York (2001).
- [4] A. BÖTTCHER, B. SILBERMANN. *Introduction to large truncated Toeplitz matrices*. Springer-Verlag, New York (1999).
- [5] J.A. COTTRELL, T.J.R. HUGHES, Y. BAZILEVS. *Isogeometric analysis: toward integration of CAD and FEA*. John Wiley & Sons (2009).
- [6] M. DONATELLI, C. GARONI, C. MANNI, S. SERRA-CAPIZZANO, H. SPELEERS. *Robust and optimal multi-iterative techniques for IgA Galerkin linear systems*. Comput. Methods Appl. Mech. Engrg. **284** (2015) 230–264.
- [7] M. DONATELLI, C. GARONI, C. MANNI, S. SERRA-CAPIZZANO, H. SPELEERS. *Robust and optimal multi-iterative techniques for IgA collocation linear systems*. Comput. Methods Appl. Mech. Engrg. **284** (2015) 1120–1146.
- [8] M. DONATELLI, C. GARONI, C. MANNI, S. SERRA-CAPIZZANO, H. SPELEERS. *Spectral analysis and spectral symbol of matrices in isogeometric collocation methods*. Math. Comp. (2015) to appear.
- [9] M. DONATELLI, C. GARONI, C. MANNI, S. SERRA-CAPIZZANO, H. SPELEERS. *Symbol-based multigrid methods for Galerkin B-spline isogeometric analysis*. Tech. Report TW650, Dept. Computer Science, KU Leuven (2014).
- [10] C. GARONI. *Structured matrices coming from PDE Approximation Theory: spectral analysis, spectral symbol and design of fast iterative solvers*. Ph.D. Thesis in Mathematics of Computation, University of Insubria, Como, Italy (2014).
- [11] C. GARONI, C. MANNI, F. PELOSI, S. SERRA-CAPIZZANO, H. SPELEERS. *On the spectrum of stiffness matrices arising from isogeometric analysis*. Numer. Math. **127** (2014) 751–799.
- [12] C. GARONI, S. SERRA-CAPIZZANO, D. SESANA. *Tools for determining the asymptotic spectral distribution of non-Hermitian perturbations of Hermitian matrix-sequences and applications*. Integr. Equ. Oper. Theory **81** (2015) 213–225.
- [13] C. GARONI, S. SERRA-CAPIZZANO, D. SESANA. *Spectral analysis and spectral symbol of  $d$ -variate  $\mathbb{Q}_p$  Lagrangian FEM stiffness matrices*. Tech. Report 2014-021, Dept. Information Technology, Uppsala University, Sweden (2014).
- [14] C. GARONI, S. SERRA-CAPIZZANO, P. VASSALOS. *A general tool for determining the asymptotic spectral distribution of Hermitian matrix-sequences*. Oper. Matrices (2015) to appear.
- [15] L. GOLINSKII, S. SERRA-CAPIZZANO. *The asymptotic properties of the spectrum of nonsymmetrically perturbed Jacobi matrix sequences*. J. Approx. Theory **144** (2007) 84–102.
- [16] L. HÖRMANDER. *Pseudo-differential operators and non-elliptic boundary problems*. Annals of Math. **83** (Second Series) (1966) 129–209.
- [17] W. RUDIN. *Real and complex analysis*. 3<sup>rd</sup> edition, McGraw-Hill, Singapore (1987).
- [18] S. SERRA-CAPIZZANO. *Distribution results on the algebra generated by Toeplitz sequences: a finite dimensional approach*. Linear Algebra Appl. **328** (2001) 121–130.
- [19] S. SERRA-CAPIZZANO. *Spectral behavior of matrix sequences and discretized boundary value problems*. Linear Algebra Appl. **337** (2001) 37–78.
- [20] S. SERRA-CAPIZZANO. *Generalized locally Toeplitz sequences: spectral analysis and applications to discretized partial differential equations*. Linear Algebra Appl. **366** (2003) 371–402.
- [21] S. SERRA-CAPIZZANO. *The GLT class as a generalized Fourier analysis and applications*. Linear Algebra Appl. **419** (2006) 180–233.

- [22] S. SERRA-CAPIZZANO, C. TABLINO POSSIO. *Analysis of preconditioning strategies for collocation linear systems*. Linear Algebra Appl. **369** (2003) 41–75.
- [23] P. TILLI. *Locally Toeplitz sequences: spectral properties and applications*. Linear Algebra Appl. **278** (1998) 91–120.
- [24] E.E. TYRTYSHNIKOV. *A unifying approach to some old and new theorems on distribution and clustering*. Linear Algebra Appl. **232** (1996) 1–43.