

# Spectral analysis of $\mathbb{P}_k$ Finite Element matrices in the case of Friedrichs–Keller triangulations via Generalized Locally Toeplitz technology

Ryma Imene Rahla<sup>1</sup> | Stefano Serra-Capizzano<sup>2,3</sup> | Cristina Tablino-Possio<sup>4</sup> 

<sup>1</sup>Ecole Nationale d'Ingénieurs de Tunis, Tunis, Tunisia

<sup>2</sup>Dipartimento di Scienze Umane e dell'Innovazione per il Territorio - INDAM Unit - Dipartimento di Scienza e Alta Tecnologia, Università dell'Insubria - Sede di Como, Como, Italy

<sup>3</sup>Division of Scientific Computing, Department of Information Technology, Uppsala University, Uppsala, Sweden

<sup>4</sup>Dipartimento di Matematica e Applicazioni, Università di Milano Bicocca, Milano, Italy

## Correspondence

Cristina Tablino-Possio, Dipartimento di Matematica e Applicazioni, Università di Milano Bicocca, Milano, Italy.  
Email: cristina.tablinopossio@unimib.it

## Summary

In the present article, we consider a class of elliptic partial differential equations with Dirichlet boundary conditions and where the operator is  $\operatorname{div}(-a(\mathbf{x})\nabla\cdot)$ , with  $a$  continuous and positive over  $\overline{\Omega}$ ,  $\Omega$  being an open and bounded subset of  $\mathbb{R}^d$ ,  $d \geq 1$ . For the numerical approximation, we consider the classical  $\mathbb{P}_k$  Finite Elements, in the case of Friedrichs–Keller triangulations, leading, as usual, to sequences of matrices of increasing size. The new results concern the spectral analysis of the resulting matrix-sequences in the direction of the global distribution in the Weyl sense, with a concise overview on localization, clustering, extremal eigenvalues, and asymptotic conditioning. We study in detail the case of constant coefficients on  $\Omega = (0, 1)^2$  and we give a brief account in the more involved case of variable coefficients and more general domains. Tools are drawn from the Toeplitz technology and from the rather new theory of Generalized Locally Toeplitz sequences. Numerical results are shown for a practical evidence of the theoretical findings.

## KEYWORDS

finite element approximations, matrix-sequences, spectral analysis

## MOS SUBJECT CLASSIFICATION

65N30; 15A18; 47B35; 15A12; 65F10

## 1 | INTRODUCTION

The article deals with the spectral analysis of matrix-sequences arising in the  $\mathbb{P}_k$  Lagrangian finite element approximation of the elliptic problem

$$\begin{cases} \operatorname{div}(-a(\mathbf{x})\nabla u) = f, & \mathbf{x} \in \Omega \subseteq \mathbb{R}^d, \\ u|_{\partial\Omega} = 0, \end{cases} \quad (1)$$

with  $\Omega$  bounded connected subset of  $\mathbb{R}^d$ ,  $d \geq 1$ , having smooth boundaries for  $d \geq 2$ , and  $a$  being continuous and positive on  $\overline{\Omega}$ . Our theoretical analysis focuses on the case of stiffness matrix-sequences  $\{A_n\}_n$  related to  $\mathbb{P}_k$  Finite element approximations on uniform structured meshes,<sup>1–4</sup> such as Friedrichs–Keller triangulations, in which context the powerful spectral tools derived from the Toeplitz theory<sup>5–9</sup> greatly facilitate the required spectral analysis.

We give a detailed analysis in the case where  $a(\mathbf{x}) \equiv 1$  and then we sketch the general setting, by considering a Riemann integrable diffusion coefficient  $a$  and/or a domain  $\Omega$  not necessarily of Cartesian structure. We recall that the same type of analysis of the linear Finite Elements in two dimensions is already considered in References 10 and 11 for the same equation considered in this note, while coupled partial differential equations (PDEs) with stable pairs of Finite Element approximations again in two dimension are considered in Reference 12. It is worth noticing the systematic work in Reference 13, where the case of tensor rectangular Finite Element approximations  $\mathbb{Q}_{\mathbf{k}}$  of any degree  $\mathbf{k} = (k_1, \dots, k_d)$  and of any dimensionality  $d \geq 1$  is studied.

Here, following the pattern indicated in Reference 13, we start a systematic approach for the Finite Element approximations  $\mathbb{P}_k$  for  $k \geq 2$  and for  $d = 2$ . The analysis for  $d = 1$  is contained in Reference 13 trivially because  $\mathbb{Q}_k \equiv \mathbb{P}_k$  for every  $k \geq 1$ , while for  $d = 2$ , and even more for  $d \geq 3$ , the situation is greatly complicated by the fact that we do not encounter a tensor structure. Nevertheless, the spectral picture is quite similar and the obtained information in terms of spectral symbol is sufficient for deducing a quite accurate analysis concerning the distribution and the extremal behavior of the eigenvalues of the resulting matrix-sequences.

More in details, regarding the resulting stiffness matrices, we will consider the following items, from the perspective of (block) multilevel Toeplitz operators<sup>5,14</sup> and (block) Generalized Locally Toeplitz (GLT) sequences:<sup>6,7</sup>

- spectral distribution in the Weyl sense,
- spectral clustering,

with a concise analysis also of the extremal eigenvalues, conditioning, spectral localization, and where the final goal is

- the analysis and the design of fast iterative solvers for the associated linear systems.

We recall that the spectral distribution and the clustering results represent key ingredients in the design and in the convergence analysis of specialized multigrid methods and preconditioned Krylov solvers<sup>15</sup> such as preconditioned conjugate gradient (PCG); see Reference 7, Subsection 3.7 and References 16–22. In fact, the knowledge of the spectral distribution is the key for explaining the superlinear convergence history of (P)CG, thus improving the classical bounds; see Reference 17 and references therein. Most of this article is actually focused on the identification of the spectral symbol via the GLT technology and hence on the first item.

## 1.1 | A comparison with the case of different approximation techniques

This subsection is dedicated to make a short technical comparison with related works, when different approximation techniques are considered. Both the similarities and the differences are described in order to provide a clear global picture.

Recently, a very close spectral analysis has been conducted for the stiffness/collocation matrices coming from the  $\mathbf{k}$ -degree B-spline Isogeometric Analysis (IgA) approximation of maximal smoothness<sup>23</sup> of (1).<sup>24,25</sup> The same kind of analysis in the case of  $\mathbb{Q}_{\mathbf{k}}$  Finite Elements approximating again (1) can be found in Reference 13. A review comparing the latter two approaches with a language tailored for engineers can be found in Reference 26.

In the IgA case, the (spectral) symbol  $f_{\text{IgA}_{\mathbf{k}}}$  describing the spectral distribution is a scalar-valued  $d$ -variate function defined over  $[-\pi, \pi]^d$ , and so the eigenvalues of the IgA discretization matrices are approximated by a uniform sampling of  $f_{\text{IgA}_{\mathbf{k}}}$  over  $[-\pi, \pi]^d$ . In this context, the surprising behavior is that, when all the spline degrees  $k$  increase,  $f_{\text{IgA}_{\mathbf{k}}}(\boldsymbol{\theta})$  collapses exponentially to zero at all points  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_d)$  with some component  $\theta_j = \pi$ . In view of the interpretation based on the theory of Toeplitz matrices and matrix algebras, this phenomenon implies that the IgA matrices are ill-conditioned, not only in the low frequencies (as expected), but also in the high frequencies, like in the approximation of integral operators.<sup>27</sup> The explicit use of this spectral information allowed the design of ad hoc iterative solvers with an optimal convergence rate, substantially independent of  $\mathbf{k}$  and  $d$ .<sup>18–20</sup>

In the  $\mathbb{Q}_{\mathbf{k}}$  Lagrangian setting, we are still able to identify the spectral distribution, as for the IgA case. The related symbol  $f_{\mathbb{Q}_{\mathbf{k}}}$  is  $d$ -variate and defined on  $[-\pi, \pi]^d$ , but the surprise is that  $f_{\mathbb{Q}_{\mathbf{k}}}$  is a  $N(\mathbf{k}) \times N(\mathbf{k})$  Hermitian matrix-valued function, with  $\mathbf{k} = (k_1, \dots, k_d)$  vector of the partial degrees in the different directions,  $N(\mathbf{k}) = \prod_{j=1}^d k_j$  (a similar situation is encountered in dealing with discontinuous Galerkin methods<sup>28,29</sup>). No specific pathologies regarding  $f_{\mathbb{Q}_{\mathbf{k}}}$  are observed for large  $\mathbf{k}$  at the points  $\boldsymbol{\theta}$  such that  $\theta_j = \pi$  for some  $j$ , implying that the source of ill-conditioning, with respect to the

fineness parameters, is only in the low frequencies. However, exactly as in the present  $\mathbb{P}_k$  setting, where  $k$  is a scalar indicating the global polynomial degree, we observe a serious problem of dimensionality, since, already for moderate  $k$  and  $d$ , the quantity  $N(k, d) = k^d$  is very large.

More specifically, the problem is that the spectrum of the  $\mathbb{Q}_k$  Lagrangian Finite Element stiffness matrices is split into  $N(\mathbf{k}) = \prod_{j=1}^d k_j$  subsets, or branches in the engineering terminology,<sup>30-33</sup> of approximately the same cardinality and the  $i$ th branch is approximately a uniform sampling of the scalar-valued function  $\lambda_i(f_{\mathbb{Q}_k})$ ,  $i = 1, \dots, N(\mathbf{k})$ . The exponential scattering (in  $\mathbf{k}$  and  $d$ ) of the eigenvalue functions  $\lambda_i(f_{\mathbb{Q}_k})$  provides an explanation of the difficulties encountered by the solvers in the literature, already for moderate  $\mathbf{k}$  and  $d$ . Indeed, it is relatively easy to design a mesh-independent solver, but the dependence on  $\mathbf{k}$  and  $d$  is generally bad. In the following we will also use the symbol  $\mathbb{Q}_k$  indicating that  $\mathbf{k} = (k, \dots, k)$ , that is,  $k_1 = k_2 = \dots = k_d = k$ : In the present case  $N(\mathbf{k}) = N(k, d) = k^d$ .

At this point it is worthwhile stressing that a cardinality of the branches equal to  $N(\mathbf{k})$  is expected in the tensor setting  $\mathbb{Q}_k$ ,  $\mathbf{k} = (k_1, \dots, k_d)$ , while it is somehow a surprise with the current choice of  $\mathbb{P}_k$  Finite Elements, where  $k$  is a scalar indicating the global degree and  $N(k, d) = k^d$ . We have in fact checked this formula only for  $d = 2$  and  $k = 2, 3, 4$  and hence a deeper analysis of this phenomenon will be the subject of future investigations.

## 1.2 | Structure and challenges of the article

In Section 2, we present the standard Galerkin approximation of (1) by  $\mathbb{P}_k$  Lagrangian Finite Elements with  $d = 2$ . Section 3 contains preliminaries concerning spectral distribution and clustering, special matrix structures such as multilevel block Toeplitz/Circulant matrices, multilevel block diagonal structures, zero-distributed sequences, and the basics of the theory of multilevel block GLT sequences. Sections 4 and 5 describe the specific 1D, 2D approximations with the related identification of the symbol for the different choice of parameters, while Section 6 is devoted to the study of the analytical features of the symbol and to the implications in terms of spectral distribution, clustering, localization, extremal eigenvalues, and conditioning. A concise account on the case of variable coefficients and non-Cartesian domains is contained in Section 7. Section 8 deals with preconditioning and complexity issues. Finally, Section 9 contains concluding remarks, open problems, and perspectives.

In nontechnical terms we report very concisely what is the challenge, how the challenge is tackled, and the relevance of this task:

- Challenge: we extend the spectral analysis to the case of  $\mathbb{P}_k$  discretizations of the diffusion equation. This discretization does not provide the tensor product structure observed in the  $\mathbb{Q}_k$  setting.
- Methodology: we use GLT technology described for nonexperts to derive the spectral distribution.
- Relevance: the final target is to develop preconditioners and specialized multigrid methods for these discretizations, especially higher-order discretizations and for variable coefficient diffusion, as done in Reference 34 in the context of multigrid procedures for  $\mathbb{Q}_k$  Finite Elements.

## 2 | FINITE ELEMENT APPROXIMATION

Problem (1) can be formulated in variational form as follows:

$$\text{find } u \in H_0^1(\Omega) \text{ such that } \int_{\Omega} (a \nabla u \cdot \nabla \varphi) = \int_{\Omega} f \varphi \quad \text{for all } \varphi \in H_0^1(\Omega), \quad (2)$$

where  $H_0^1(\Omega)$  is the space of square integrable functions vanishing on  $\partial\Omega$ , with square integrable weak derivatives. We assume that  $\Omega \subseteq \mathbb{R}^2$  is a bounded connected set with smooth boundaries (in practice in our numerical tests  $\Omega$  will be a polygonal domain). Furthermore, we make the following hypotheses on the coefficients:

$$a \in C(\overline{\Omega}), \text{ with } a(\mathbf{x}) \geq a_0 > 0, \quad \text{and } f \in L^2(\Omega), \quad (3)$$

so that existence and uniqueness for problem (2) is guaranteed. Hereafter, we consider  $\mathbb{P}_k$  Lagrangian Finite Element approximation of problem (2). To this end, let  $\mathcal{T}_h = \{K\}$  be a usual Finite Element partition of  $\Omega$  into triangles, with

$h_K = \text{diam}(K)$  and  $h = \max_K h_K$ , and let  $V_h \subset H_0^1(\Omega)$  be the space of  $\mathbb{P}_k$  Lagrangian Finite Element, that is,

$$V_h = \{\varphi_h : \bar{\Omega} \rightarrow \mathbb{R} \text{ s.t. } \varphi_h \text{ is continuous, } \varphi_h|_K \text{ is a polynomial of degree less or equal to } k, \text{ and } \varphi_h|_{\partial\Omega} = 0\}.$$

The Finite Element approximation of problem (2) reads as follows:

$$\text{find } u_h \in V_h \text{ such that } \int_{\Omega} (a \nabla u_h \cdot \nabla \varphi_h) = \int_{\Omega} f \varphi_h \quad \text{for all } \varphi_h \in V_h. \quad (4)$$

For each internal node  $i$  of the mesh  $\mathcal{T}_h$ , meaning both vertices and additional nodal values associated with the  $\mathbb{P}_k$  approximation, let  $\varphi_i \in V_h$  be such that  $\varphi_i(\text{node } i) = 1$ , and  $\varphi_i(\text{node } j) = 0$  if  $i \neq j$ . Then, the collection of all  $\varphi_i$ 's is a basis for  $V_h$  and we denote by  $n(h)$  its dimension. Then, we write  $u_h$  as  $u_h = \sum_{j=1}^{n(h)} u_j \varphi_j$  and the variational equation (4) becomes an algebraic linear system:

$$\sum_{j=1}^{n(h)} \left( \int_{\Omega} a \nabla \varphi_j \cdot \nabla \varphi_i \right) u_j = \int_{\Omega} f \varphi_i, \quad i = 1, \dots, n(h). \quad (5)$$

The aim of this article is to analyze the spectral properties of the matrix-sequences  $\{A_n(a, \Omega, \mathbb{P}_k)\}_n$  arising in the quoted linear systems (5), both from the theoretical and numerical point of view.

### 3 | PRELIMINARIES ON MATRIX-SEQUENCES: TOEPLITZ AND GLT STRUCTURES

As recalled in the introduction, the main target of the article is the spectral analysis of the considered stiffness matrices, from the perspective of (block) multilevel Toeplitz operators and (block) GLT sequences, with special attention to

- spectral distribution in the Weyl sense,
- spectral clustering,

and with a concise analysis also of the extremal eigenvalues, conditioning, and spectral localization.

In this section, which is divided into three subsections, we furnish the tools for handling such items. In the first subsection we introduce the notion of clustering and distribution for general matrix-sequences (Subsection 3.1). The other two are devoted to multilevel block Toeplitz matrices and further special matrix-sequences (Subsection 3.2) and to the  $*$ -algebra of GLT sequences (Subsection 3.3).

#### 3.1 | Clustering and distribution

This subsection is devoted to the notion of distribution (with a matrix-valued symbol), both in the sense of the eigenvalues and singular values. The notion of clustering, both in the sense of eigenvalues and singular values, can be seen as a special case of the distribution notions.

• *Sequences of matrices and block matrix-sequences.* Throughout this article, a sequence of matrices is any sequence of the form  $\{A_n\}_n$ , where  $A_n$  is a square matrix of size  $d_n$  and  $d_n \rightarrow \infty$  as  $n \rightarrow \infty$ . Let  $r \geq 1$  be a fixed positive integer independent of  $n$ ; an  $r$ -block matrix-sequence (or simply a matrix-sequence if  $r$  can be inferred from the context or we do not need/want to specify it) is a special sequence of matrices  $\{A_n\}_n$  in which the size of  $A_n$  is  $d_n = r\phi_n$ , with  $\{\phi_n\}_n$  being a sequence of positive integers.

• *Singular value and eigenvalue distribution of a sequence of matrices.* Let  $\mu_t$  be the Lebesgue measure in  $\mathbb{R}^t$ ,  $t \geq 1$ . Throughout this article, all the terminology from measure theory (such as “measurable set,” “measurable function,” “a.e.,” etc.) is referred to as the Lebesgue measure. A matrix-valued function  $f : D \subseteq \mathbb{R}^t \rightarrow \mathbb{C}^{r \times r}$  is said to be measurable (resp., continuous, Riemann-integrable, in  $L^p(D)$ , etc.) if its components  $f_{\alpha\beta} : D \rightarrow \mathbb{C}$ ,  $\alpha, \beta = 1, \dots, r$ , are measurable (resp., continuous, Riemann-integrable, in  $L^p(D)$ , etc.). We denote by  $C_c(\mathbb{R})$  (resp.,  $C_c(\mathbb{C})$ ) the space of continuous

complex-valued functions with bounded support defined on  $\mathbb{R}$  (resp.,  $\mathbb{C}$ ). If  $A \in \mathbb{C}^{m \times m}$ , the singular values and the eigenvalues of  $A$  are denoted by  $\sigma_1(A), \dots, \sigma_m(A)$  and  $\lambda_1(A), \dots, \lambda_m(A)$ , respectively.

**Definition 1.** Let  $\{A_n\}_n$  be a sequence of matrices, with  $A_n$  of size  $d_n$ , and let  $f : D \subset \mathbb{R}^t \rightarrow \mathbb{C}^{r \times r}$  be a measurable function defined on a set  $D$  with  $0 < \mu_t(D) < \infty$ .

- We say that  $\{A_n\}_n$  has a (asymptotic) singular value distribution described by  $f$ , and we write  $\{A_n\}_n \sim_\sigma f$ , if

$$\lim_{n \rightarrow \infty} \frac{1}{d_n} \sum_{i=1}^{d_n} F(\sigma_i(A_n)) = \frac{1}{\mu_t(D)} \int_D \frac{\sum_{i=1}^r F(\sigma_i(f(\mathbf{x})))}{r} d\mathbf{x}, \quad \forall F \in C_c(\mathbb{R}). \quad (6)$$

- We say that  $\{A_n\}_n$  has a (asymptotic) spectral (or eigenvalue) distribution described by  $f$ , and we write  $\{A_n\}_n \sim_\lambda f$ , if

$$\lim_{n \rightarrow \infty} \frac{1}{d_n} \sum_{i=1}^{d_n} F(\lambda_i(A_n)) = \frac{1}{\mu_t(D)} \int_D \frac{\sum_{i=1}^r F(\lambda_i(f(\mathbf{x})))}{r} d\mathbf{x}, \quad \forall F \in C_c(\mathbb{C}). \quad (7)$$

If  $\{A_n\}_n$  has both a singular value and an eigenvalue distribution described by  $f$ , we write  $\{A_n\}_n \sim_{\sigma, \lambda} f$ .

We note that Definition 1 is well-posed because the functions

$$\mathbf{x} \mapsto \sum_{i=1}^r F(\sigma_i(f(\mathbf{x}))) \quad \text{and} \quad \mathbf{x} \mapsto \sum_{i=1}^r F(\lambda_i(f(\mathbf{x})))$$

are measurable. Whenever we write a relation such as  $\{A_n\}_n \sim_\sigma f$  or  $\{A_n\}_n \sim_\lambda f$ , it is understood that  $f$  is as in Definition 1, that is,  $f$  is a measurable function defined on a subset  $D$  of some  $\mathbb{R}^t$  with  $0 < \mu_t(D) < \infty$  and taking values in  $\mathbb{C}^{r \times r}$  for some  $r \geq 1$ . The informal meaning behind the spectral distribution (7) is the following: If  $f$  is continuous, then a suitable ordering of the eigenvalues  $\{\lambda_j(A_n)\}_{j=1, \dots, d_n}$ , assigned in correspondence with an equispaced grid on  $D$ , reconstructs approximately the  $r$  surfaces  $\mathbf{x} \mapsto \lambda_i(f(\mathbf{x})), i = 1, \dots, r$ . For instance, if  $t = 1$ ,  $d_n = nr$ , and  $D = [a, b]$ , then the eigenvalues of  $A_n$  are approximately equal to  $\lambda_i(f(a + j(b - a)/n)), j = 1, \dots, n, i = 1, \dots, r$ ; if  $t = 2$ ,  $d_n = n^2 r$ , and  $D = [a_1, b_1] \times [a_2, b_2]$ , then the eigenvalues of  $A_n$  are approximately equal to  $\lambda_i(f(a_1 + j_1(b_1 - a_1)/n, a_2 + j_2(b_2 - a_2)/n)), j_1, j_2 = 1, \dots, n, i = 1, \dots, r$  (and so on for  $t \geq 3$ ). This type of information is useful in engineering applications,<sup>26</sup> for example, for the computation of the relevant vibrations, and in the analysis of the (asymptotic) convergence speed of iterative solvers for large linear systems or for improving the convergence rate by, for example, the design of appropriate preconditioners.<sup>17,10</sup>

The next theorem gives useful tools for computing the spectral distribution of sequences formed by Hermitian matrices. For the related proof, we refer the reader to Reference 22, Theorem 4.3. In what follows, the conjugate transpose of the matrix  $A$  is denoted by  $A^*$ . If  $A \in \mathbb{C}^{m \times m}$  and  $1 \leq p \leq \infty$ , we denote by  $\|A\|_p$  the Schatten  $p$ -norm of  $A$ , that is, the  $p$ -norm of the vector  $(\sigma_1(A), \dots, \sigma_m(A))$ . The Schatten  $\infty$ -norm  $\|A\|_\infty$  is the largest singular value of  $A$  and coincides with the spectral norm  $\|A\|$ . The Schatten 1-norm  $\|A\|_1$  is the sum of the singular values of  $A$  and is often referred to as the trace-norm of  $A$ . The Schatten 2-norm  $\|A\|_2$  coincides with the Frobenius norm of  $A$ . For more on Schatten  $p$ -norms, see Reference 35.

**Theorem 1.** Let  $\{X_n\}_n$  be a sequence of matrices, with  $X_n$  Hermitian of size  $d_n$ , and let  $\{P_n\}_n$  be a sequence such that  $P_n \in \mathbb{C}^{d_n \times \delta_n}$ ,  $P_n^* P_n = I_{\delta_n}$ ,  $\delta_n \leq d_n$  and  $\delta_n/d_n \rightarrow 1$  as  $n \rightarrow \infty$ . Then,  $\{X_n\}_n \sim_{\sigma, \lambda} \kappa$  if and only if  $\{P_n^* X_n P_n\}_n \sim_{\sigma, \lambda} \kappa$ .

Now we turn to the definition of clustering. For  $z \in \mathbb{C}$  and  $\epsilon > 0$ , let  $B(z, \epsilon)$  the disk with center  $z$  and radius  $\epsilon$ ,  $B(z, \epsilon) := \{w \in \mathbb{C} : |w - z| < \epsilon\}$ . For  $S \subseteq \mathbb{C}$  and  $\epsilon > 0$ , we denote by  $B(S, \epsilon)$  the  $\epsilon$ -expansion of  $S$ , defined as  $B(S, \epsilon) := \cup_{z \in S} B(z, \epsilon)$ .

**Definition 2.** Let  $\{X_n\}_n$  be a sequence of matrices, with  $X_n$  of size  $d_n$  tending to infinity, and let  $S \subseteq \mathbb{C}$  be a nonempty closed subset of  $\mathbb{C}$ .  $\{X_n\}_n$  is *strongly clustered* at  $S$  in the sense of the eigenvalues if, for each  $\epsilon > 0$ , the number of eigenvalues of  $X_n$  outside  $B(S, \epsilon)$  is bounded by a constant  $q_\epsilon$  independent of  $n$ . In symbols,

$$q_\epsilon(n, S) := \#\{j \in \{1, \dots, d_n\} : \lambda_j(X_n) \notin B(S, \epsilon)\} = O(1), \quad \text{as } n \rightarrow \infty.$$

$\{X_n\}_n$  is *weakly clustered* at  $S$  if, for each  $\epsilon > 0$ ,

$$q_\epsilon(n, S) = o(d_n), \quad \text{as } n \rightarrow \infty.$$

If  $\{X_n\}_n$  is strongly or weakly clustered at  $S$  and  $S$  is not connected, then the connected components of  $S$  are called subclusters.

Recall that, for a measurable function  $g : D \subseteq \mathbb{R}^t \rightarrow \mathbb{C}$ , the essential range of  $g$  is defined as  $\mathcal{ER}(g) := \{z \in \mathbb{C} : \mu_t(\{g \in B(z, \epsilon)\}) > 0 \text{ for all } \epsilon > 0\}$ , where  $\{g \in B(z, \epsilon)\} := \{x \in D : g(x) \in B(z, \epsilon)\}$ .  $\mathcal{ER}(g)$  is always closed; moreover, if  $g$  is continuous and  $D$  is contained in the closure of its interior, then  $\mathcal{ER}(g)$  coincides with the closure of the image of  $g$ .

Now, if  $\{X_n\}_n \sim_\lambda f$  (with  $\{X_n\}_n, f$  as in Definition 1), then, by Reference 36, Theorem 4.2,  $\{X_n\}_n$  is weakly clustered at the essential range of  $f$ , defined as the union of the essential ranges of the eigenvalue functions  $\lambda_i(f)$ ,  $i = 1, \dots, r$ :  $\mathcal{ER}(f) := \cup_{i=1}^r \mathcal{ER}(\lambda_i(f))$ .

### 3.2 | Multilevel block Toeplitz/Circulant/diagonal matrices and zero-distributed matrix-sequences

In this subsection we introduce three types of matrix structures. The first two have an algebraic definition for every fixed dimension (multilevel block Toeplitz/Circulant matrices and multilevel block diagonal structures), while the last has only an asymptotic sense (zero-distributed matrix-sequences). In any case, we will be interested in matrix-sequences consisting of these three matrix structures, especially when defining the basics of the theory of multilevel block GLT sequences.

• *Multilevel block Toeplitz/Circulant matrices.* We first briefly summarize the definition and relevant properties of multilevel block Toeplitz matrices we will face in the following sections.

Given  $\mathbf{n} \in \mathbb{N}^d$ , a matrix of the form

$$[A_{\mathbf{i}-\mathbf{j}}]_{\mathbf{i}, \mathbf{j}=\mathbf{e}}^{\mathbf{n}} \in \mathbb{C}^{N(\mathbf{n})r \times N(\mathbf{n})r}$$

with  $\mathbf{e}$  vector of all ones, with blocks  $A_{\mathbf{k}} \in \mathbb{C}^{r \times r}$ ,  $\mathbf{k} = -(\mathbf{n} - \mathbf{e}), \dots, \mathbf{n} - \mathbf{e}$ , is called a multilevel block Toeplitz matrix, or, more precisely, a  $d$ -level  $r$ -block Toeplitz matrix. Given a matrix-valued function  $f : [-\pi, \pi]^d \rightarrow \mathbb{C}^{r \times r}$  in  $L^1([-\pi, \pi]^d)$ , we denote its Fourier coefficients by

$$\hat{f}_{\mathbf{k}} = \frac{1}{(2\pi)^d} \int_{[-\pi, \pi]^d} f(\boldsymbol{\theta}) e^{-i(\mathbf{k}, \boldsymbol{\theta})} d\boldsymbol{\theta} \in \mathbb{C}^{r \times r}, \quad \mathbf{k} \in \mathbb{Z}^d, \quad (8)$$

where the integrals are computed componentwise and  $(\mathbf{k}, \boldsymbol{\theta}) = k_1\theta_1 + \dots + k_d\theta_d$ . For every  $\mathbf{n} \in \mathbb{N}^d$ , the  $\mathbf{n}$ th Toeplitz matrix associated with  $f$  is defined as

$$T_{\mathbf{n}}(f) := [\hat{f}_{\mathbf{i}-\mathbf{j}}]_{\mathbf{i}, \mathbf{j}=\mathbf{e}}^{\mathbf{n}} \quad (9)$$

or, equivalently, as

$$T_{\mathbf{n}}(f) = \sum_{|\mathbf{j}_1| < n_1} \dots \sum_{|\mathbf{j}_d| < n_d} [J_{n_1}^{(\mathbf{j}_1)} \otimes \dots \otimes J_{n_d}^{(\mathbf{j}_d)}] \otimes \hat{f}_{(\mathbf{j}_1, \dots, \mathbf{j}_d)} \quad (10)$$

where  $\otimes$  denotes the (Kronecker) tensor product of matrices, while  $J_m^{(l)}$  is the matrix of order  $m$  whose  $(i, j)$  entry equals 1 if  $i - j = l$  and zero otherwise. We call  $\{T_{\mathbf{n}}(f)\}_{\mathbf{n} \in \mathbb{N}^d}$  the family of (multilevel block) Toeplitz matrices associated with  $f$ , which, in turn, is called the generating function of  $\{T_{\mathbf{n}}(f)\}_{\mathbf{n} \in \mathbb{N}^d}$ . In perfect analogy we define multilevel block Circulant matrices. Given  $\mathbf{n} \in \mathbb{N}^d$ , a matrix of the form

$$[A_{(\mathbf{i}-\mathbf{j}) \bmod \mathbf{n}}]_{\mathbf{i}, \mathbf{j}=\mathbf{e}}^{\mathbf{n}} \in \mathbb{C}^{N(\mathbf{n})r \times N(\mathbf{n})r}$$



with  $\mathbf{e}$  vector of all ones, with blocks  $A_{\mathbf{k}} \in \mathbb{C}^{r \times r}$ ,  $\mathbf{k} = \mathbf{0}, \dots, \mathbf{n} - \mathbf{e}$ , is called a multilevel block Circulant matrix, or, more precisely, a  $d$ -level  $r$ -block Circulant matrix.

The  $\mathbf{n}$ th Circulant matrix associated with  $f$  is defined as

$$C_{\mathbf{n}}(f) = \sum_{|j_1| < n_1} \cdots \sum_{|j_d| < n_d} [Z_{n_1}^{(j_1)} \otimes \cdots \otimes Z_{n_d}^{(j_d)}] \otimes \hat{f}_{(j_1, \dots, j_d)}, \quad (11)$$

where  $Z_m^{(l)} = [Z_m^{(1)}]^l$  is the matrix of order  $m$  whose  $(i, j)$  entry equals 1 if  $(i - j) \bmod m = l$  and zero otherwise.

• **Block diagonal sampling matrices.** For  $n \in \mathbb{N}$ ,  $d = 1$ , and  $a : [0, 1] \rightarrow \mathbb{C}^{r \times r}$ , we define the block diagonal sampling matrix  $D_n(a)$  as the block diagonal matrix

$$D_n(a) = \text{diag}_{i=1, \dots, n} a\left(\frac{i}{n}\right) = \begin{bmatrix} a\left(\frac{1}{n}\right) & & \\ & a\left(\frac{2}{n}\right) & \\ & & \ddots \\ & & & a(1) \end{bmatrix} \in \mathbb{C}^{rn \times rn}.$$

For a general dimensionality  $d \geq 2$ , we consider  $a : [0, 1]^d \rightarrow \mathbb{C}^{r \times r}$ ,  $\mathbf{n} = (n_1, \dots, n_d)$  and we define the block multilevel diagonal sampling matrix  $D_{\mathbf{n}}(a)$  as the block diagonal matrix

$$D_{\mathbf{n}}(a) = \text{diag}_{\mathbf{i}=\mathbf{e}, \dots, \mathbf{n}} a\left(\frac{\mathbf{i}}{\mathbf{n}}\right) \in \mathbb{C}^{rN(\mathbf{n}) \times rN(\mathbf{n})},$$

where the multiindex  $\mathbf{i}/\mathbf{n}$  has to be intended as  $(i_1/n_1, \dots, i_d/n_d)$  and the ordering is the lexicographical one as in the work by E. Tyrtyshnikov (see for instance Reference 37).

• **Zero-distributed sequences.** A sequence of matrices  $\{Z_n\}_n$  such that  $\{Z_n\}_n \sim_{\sigma} 0$  is referred to as a zero-distributed sequence. Note that, for any  $r \geq 1$ ,  $\{Z_n\}_n \sim_{\sigma} 0$  is equivalent to  $\{Z_n\}_n \sim_{\sigma} O_r$  (throughout this article,  $O_m$  and  $I_m$  denote the  $m \times m$  zero matrix and the  $m \times m$  identity matrix, respectively). Proposition 1 provides an important characterization of zero-distributed sequences together with a useful sufficient condition for detecting such sequences. Throughout this article, we use the natural convention  $1/\infty = 0$ .

**Proposition 1.** Let  $\{Z_n\}_n$  be a sequence of matrices, with  $Z_n$  of size  $d_n$ , and let  $\|\cdot\|$  be the spectral norm. Then

- $\{Z_n\}_n$  is zero-distributed if and only if  $Z_n = R_n + N_n$  with  $\text{rank}(R_n)/d_n \rightarrow 0$  and  $\|N_n\| \rightarrow 0$  as  $n \rightarrow \infty$ .
- $\{Z_n\}_n$  is zero-distributed if there exists a  $p \in [1, \infty]$  such that  $\|Z_n\|_p/(d_n)^{1/p} \rightarrow 0$  as  $n \rightarrow \infty$ .

### 3.3 | The $*$ -algebra of multilevel block GLT sequences

Let  $r \geq 1$  be a fixed positive integer. An  $r$ -block GLT sequence (or simply a GLT sequence if  $r$  can be inferred from the context or we do not need/want to specify it) is a special  $r$ -block matrix-sequence  $\{A_n\}_n$  equipped with a measurable function  $\kappa : [0, 1]^d \times [-\pi, \pi]^d \rightarrow \mathbb{C}^{r \times r}$ ,  $d \geq 1$ , the so-called symbol. We use the notation  $\{A_n\}_n \sim_{\text{GLT}} \kappa$  to indicate that  $\{A_n\}_n$  is a GLT sequence with symbol  $\kappa$ . The symbol of a GLT sequence is unique in the sense that if  $\{A_n\}_n \sim_{\text{GLT}} \kappa$  and  $\{A_n\}_n \sim_{\text{GLT}} \varsigma$  then  $\kappa = \varsigma$  a.e. in  $[0, 1]^d \times [-\pi, \pi]^d$ .

The main properties of  $r$ -block GLT sequences proved in References 38 and 39 are listed below: They represent a complete characterization of GLT sequences, equivalent to the full constructive definition.

If  $A$  is a matrix, we denote by  $A^\dagger$  the Moore–Penrose pseudoinverse of  $A$  (recall that  $A^\dagger = A^{-1}$  whenever  $A$  is invertible). If  $f_m, f : D \subseteq \mathbb{R}^t \rightarrow \mathbb{C}^{r \times r}$  are measurable matrix-valued functions, we say that  $f_m$  converges to  $f$  in measure (resp., a.e., in  $L^p(D)$ , etc.) if  $(f_m)_{\alpha\beta}$  converges to  $f_{\alpha\beta}$  in measure (resp., a.e., in  $L^p(D)$ , etc.) for all  $\alpha, \beta = 1, \dots, r$ .

**GLT1.** If  $\{A_n\}_n \sim_{\text{GLT}} \kappa$  then  $\{A_n\}_n \sim_{\sigma} \kappa$ . If moreover each  $A_n$  is Hermitian then  $\{A_n\}_n \sim_{\lambda} \kappa$ .

**GLT2.** We have:

- $\{T_n(f)\}_n \sim_{\text{GLT}} \kappa(\mathbf{x}, \boldsymbol{\theta}) = f(\boldsymbol{\theta})$  if  $f : [-\pi, \pi]^d \rightarrow \mathbb{C}^{r \times r}$  is in  $L^1([-\pi, \pi]^d)$ ;
- $\{D_n(a)\}_n \sim_{\text{GLT}} \kappa(\mathbf{x}, \boldsymbol{\theta}) = a(\mathbf{x})$  if  $a : [0, 1]^d \rightarrow \mathbb{C}^{r \times r}$  is Riemann-integrable;
- $\{Z_n\}_n \sim_{\text{GLT}} \kappa(\mathbf{x}, \boldsymbol{\theta}) = O_r$  if and only if  $\{Z_n\}_n \sim_{\sigma} 0$ .

**GLT3.** If  $\{A_n\}_n \sim_{\text{GLT}} \kappa$  and  $\{B_n\}_n \sim_{\text{GLT}} \varsigma$  then:

- $\{A_n^*\}_n \sim_{\text{GLT}} \kappa^*$ ;
- $\{\alpha A_n + \beta B_n\}_n \sim_{\text{GLT}} \alpha \kappa + \beta \varsigma$  for all  $\alpha, \beta \in \mathbb{C}$ ;
- $\{A_n B_n\}_n \sim_{\text{GLT}} \kappa \varsigma$ ;
- $\{A_n^\dagger\}_n \sim_{\text{GLT}} \kappa^{-1}$  provided that  $\kappa$  is invertible a.e.

**GLT4.**  $\{A_n\}_n \sim_{\text{GLT}} \kappa$  if and only if there exist  $r$ -block GLT sequences  $\{B_{n,m}\}_n \sim_{\text{GLT}} \kappa_m$  such that  $\{B_{n,m}\}_n \xrightarrow{\text{a.c.s.}} \{A_n\}_n$  and  $\kappa_m \rightarrow \kappa$  in measure, where the a.c.s. convergence is studied in Reference 8.

Regarding notations on sequences associated with a multiindex, all the previous definitions and notions apply with the convention that  $\mathbf{n} \rightarrow \infty$  has the meaning that  $\min_j n_j \rightarrow \infty$ . Finally, it is worth noticing that in the derivations in the following sections and in all the numerical experiments, the  $d$ -index is simplified that is  $\mathbf{n} = n \cdot \mathbf{e}$ , that is,  $n_1 = n_2 = \dots = n_d = n$ . In that setting we use the simplified notation  $T_{\mathbf{n}}(f) = T_n(f)$  and  $D_{\mathbf{n}}(a) = D_n(a)$ , where the number of levels  $d$  is understood by looking at the number of variables characterizing the definition domains either of  $f$  or of  $a$ , or by simply looking at the size  $rN(n, d)$ ,  $N(n, d) = n^d$ , of the considered matrices.

#### 4 | A FEW REMARKS ON THE MONODIMENSIONAL CASE: $\mathbb{Q}_k \equiv \mathbb{P}_k$ , $d = 1$

We report some results derived in Reference 13 for the Lagrangian Finite Elements  $\mathbb{Q}_k \equiv \mathbb{P}_k$ ,  $d = 1$ . Let us consider the Lagrange polynomials  $L_0, \dots, L_k$  associated with the reference knots  $t_j = j/k, j = 0, \dots, k$ :

$$L_i(t) = \prod_{\substack{j=0 \\ j \neq i}}^k \frac{t - t_j}{t_i - t_j} = \prod_{\substack{j=0 \\ j \neq i}}^k \frac{kt - j}{i - j}, \quad i = 0, \dots, k,$$

$$L_i(t_j) = \delta_{ij}, \quad i, j = 0, \dots, k, \quad (12)$$

and let the symbol  $\langle \cdot, \cdot \rangle$  denote the scalar product in  $L^2([0, 1])$ , that is,  $\langle \varphi, \psi \rangle := \int_0^1 \varphi \psi$ . In the case  $a(x) \equiv 1$  and  $\Omega = (0, 1)$ , the  $\mathbb{Q}_k$  matrix  $A_n(a, \Omega, \mathbb{Q}_k)$  equals the matrix  $K_n^{(k)}$  in Theorem 2.

**Theorem 2.** Let  $k, n \geq 1$ . Then

$$K_n^{(k)} = \begin{bmatrix} K_0 & K_1^T & & \\ K_1 & \ddots & \ddots & \\ & \ddots & \ddots & K_1^T \\ & & K_1 & K_0 \end{bmatrix} \quad (13)$$

where the subscript “ $-$ ” means that the last row and column of the matrix in square brackets are deleted, while  $K_0, K_1$  are  $k \times k$  blocks given by

$$K_0 = \begin{bmatrix} \langle L'_1, L'_1 \rangle & \dots & \langle L'_{k-1}, L'_1 \rangle & \langle L'_k, L'_1 \rangle \\ \vdots & & \vdots & \vdots \\ \langle L'_1, L'_{k-1} \rangle & \dots & \langle L'_{k-1}, L'_{k-1} \rangle & \langle L'_k, L'_{k-1} \rangle \\ \langle L'_1, L'_k \rangle & \dots & \langle L'_{k-1}, L'_k \rangle & \langle L'_k, L'_k \rangle + \langle L'_0, L'_0 \rangle \end{bmatrix},$$

$$K_1 = \begin{bmatrix} 0 & 0 & \dots & 0 & \langle L'_0, L'_1 \rangle \\ 0 & 0 & \dots & 0 & \langle L'_0, L'_2 \rangle \\ \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & \dots & 0 & \langle L'_0, L'_k \rangle \end{bmatrix}, \quad (14)$$

with  $L_0, \dots, L_k$  being the Lagrange polynomials in (12). In particular,  $K_n^{(k)}$  is the  $(nk - 1) \times (nk - 1)$  leading principal submatrix of the block Toeplitz matrices  $T_n(f_k)$  and  $f_k : [-\pi, \pi] \rightarrow \mathbb{C}^{k \times k}$  is an Hermitian matrix-valued function given by

$$f_k(\theta) := K_0 + K_1 e^{i\theta} + K_1^T e^{-i\theta}. \quad (15)$$



An interesting property of the Hermitian matrix-valued functions  $f_k(\theta)$  defined in (15) is reported in the theorem below. From the point of view of the spectral distribution, the message is that, independently of the parameter  $k$ , the spectral symbol is of the same character as  $2 - 2 \cos(\theta)$ , which is the symbol of the basic linear Finite Elements and the most standard Finite Differences approximation.

**Theorem 3.** *Let  $k \geq 1$ , then*

$$\det(f_k(\theta)) = d_k(2 - 2 \cos(\theta)), \quad (16)$$

where  $d_k = \det([\langle L'_j, L'_i \rangle]_{i,j=1}^k) = \det([\langle L'_j, L'_i \rangle]_{i,j=1}^{k-1}) > 0$  (with  $d_1 = 1$  being the determinant of the empty matrix by convention) and  $L_0, \dots, L_k$  are the Lagrange polynomials (12).

## 5 | TWO-DIMENSIONAL CASE: $\mathbb{P}_k$ , $d = 2$ -SYMBOL DEFINITION

Hereafter, we focus on  $\mathbb{P}_k$  Lagrangian Finite Elements in the case of Friedrichs–Keller triangulations  $\{\mathcal{T}_K\}$  of the domain  $\Omega$  as reported in Figure 1. Nodes, that is both vertices and additional nodal values associated with the chosen  $\mathbb{P}_k$  approximation, are ordered in standard lexicographical way from left to right, from bottom to top.

The stiffness matrix is built by considering the standard assembling procedure with respect to the reference element  $\hat{K}$  in Figure 2. Let  $G$  be the affine transformation mapping  $\hat{K}$  onto a generic  $K \in \mathcal{T}_K$  defined as

$$G \left( \begin{bmatrix} \hat{x} \\ \hat{y} \end{bmatrix} \right) = \begin{bmatrix} (e_3)_1 & -(e_2)_1 \\ (e_3)_2 & -(e_2)_2 \end{bmatrix} \begin{bmatrix} \hat{x} \\ \hat{y} \end{bmatrix} + \begin{bmatrix} x^{v_1} \\ y^{v_1} \end{bmatrix},$$

where  $e_1 = [x^{v_3} - x^{v_2}, y^{v_3} - y^{v_2}]^T$ ,  $e_2 = [x^{v_1} - x^{v_3}, y^{v_1} - y^{v_3}]^T$ ,  $e_3 = [x^{v_2} - x^{v_1}, y^{v_2} - y^{v_1}]^T$  represent the oriented edge vectors and  $(x^{v_i}, y^{v_i})$  are the coordinates of the  $i$ th vertex  $v_i$ . Thus,

$$A_K^{El} = \left[ \int_K \nabla \varphi_j \cdot \nabla \varphi_i \right]_{ij}$$

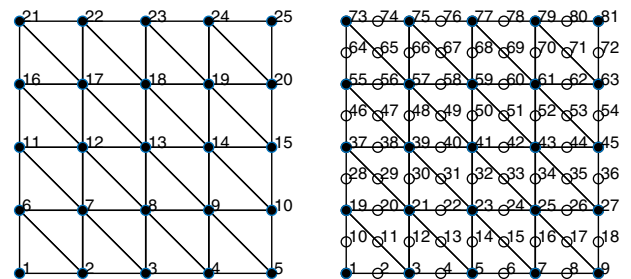
with

$$\int_K \nabla \varphi_j \cdot \nabla \varphi_i = \det(J_G(\hat{x}, \hat{y})) \int_{\hat{K}} [J_{G^{-1}}^T \hat{\nabla} \hat{\varphi}_j(\hat{x}, \hat{y})] \cdot [J_{G^{-1}}^T \hat{\nabla} \hat{\varphi}_i(\hat{x}, \hat{y})] d\hat{x} d\hat{y}, \quad (17)$$

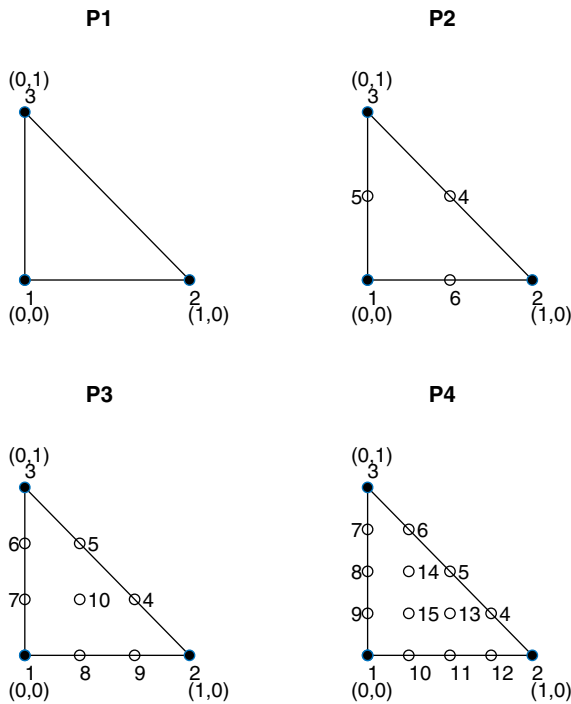
where the  $\hat{\varphi}_s$ 's are the shape functions on  $\hat{K}$ ,  $\det(J_G(\hat{x}, \hat{y})) = 2|K|$  and  $J_{G^{-1}}^T$  is the transpose of the Jacobian matrix of the inverse mapping  $G^{-1}$ , that is,

$$J_{G^{-1}}^T = \frac{1}{2|K|} \begin{bmatrix} -(e_2)_2 & -(e_3)_2 \\ (e_2)_1 & (e_3)_1 \end{bmatrix}.$$

In the present section, we will preliminarily consider the case  $a \equiv 1$  and  $\Omega = (0, 1)^2$ .



**FIGURE 1** Friedrichs–Keller meshes for  $\mathbb{P}_k$ ,  $k = 1, 2$



**FIGURE 2** Reference element  $\hat{K}$  and nodal points for  $\mathbb{P}_k$ ,  $k = 1, \dots, 4$

### 5.1 | Case $k = 1$

Even if well known, we start by considering the case  $k = 1$ , that is the one of a linear Lagrangian FE approximation. The shape functions on  $\hat{K}$  are defined as

$$\begin{aligned}\hat{\varphi}_1(\hat{x}, \hat{y}) &= -\hat{x} - \hat{y} + 1, \\ \hat{\varphi}_2(\hat{x}, \hat{y}) &= \hat{x}, \\ \hat{\varphi}_3(\hat{x}, \hat{y}) &= \hat{y},\end{aligned}\tag{18}$$

so that, according to (17), the elemental matrix for a generic triangle of the Friedrichs–Keller triangulation, that is, a right-angle triangle of constant edge  $h$ , equals

$$A_{K_1}^{El} = \frac{1}{2} \begin{bmatrix} 2 & -1 & -1 \\ -1 & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix} \quad \text{or} \quad A_{K_2}^{El} = \frac{1}{2} \begin{bmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{bmatrix},\tag{19}$$

for triangles of Type 1 (right angle in vertex 1) or Type 2 (right angle in vertex 2), respectively.

The stiffness matrix  $A_n = A_n(1, \Omega, \mathbb{P}_1)$ , that is, with  $a \equiv 1$ , is the two-level Toeplitz matrix generated by the symbol  $f_{\mathbb{P}_1}(\theta_1, \theta_2) = 4 - 2 \cos(\theta_1) - 2 \cos(\theta_2)$ . In fact,  $A_n$  is block tridiagonal, that is,

$$A_n = \text{tridiag}(A_1, A_0, A_{-1}),$$

where the triangular blocks are such that  $A_0 = \text{tridiag}(a_1^0, a_0^0, a_{-1}^0) = \text{tridiag}(-1, 4, -1)$ ,  $A_1 = A_{-1} = \text{diag}(a_1^1) = -I$ ,  $I$  being the identity matrix. Thus, we can easily read the corresponding symbol as follows

$$f_{\mathbb{P}_1}(\theta_1, \theta_2) = f_{A_0}(\theta_1) + f_{A_{-1}}(\theta_1)e^{-i\theta_2} + f_{A_1}(\theta_1)e^{i\theta_2},\tag{20}$$

with  $f_{A_0}(\theta_1) = a_0^0 + a_{-1}^0 e^{-i\theta_1} + a_1^0 e^{i\theta_1}$  and  $f_{A_1}(\theta_1) = f_{A_{-1}}(\theta_1) = a_1^1$ .

Clearly, the natural arising question is which properties are preserved in considering Lagrangian FE of higher order?

## 5.2 | Case $k = 2$

Hereafter, we will consider in full detail the case of quadratic Lagrangian FE ( $k = 2$ ), the aim being to introduce a suitable notation making easier the analysis of higher order approximations as well. By referring to the reference element (see Figure 2), we have the following shape functions

$$\begin{aligned}\hat{\phi}_1(\hat{x}, \hat{y}) &= 2\hat{x}^2 + 2\hat{y}^2 + 4\hat{x}\hat{y} - 3\hat{x} - 3\hat{y} + 1, \\ \hat{\phi}_2(\hat{x}, \hat{y}) &= \hat{x}(2\hat{x} - 1), \\ \hat{\phi}_3(\hat{x}, \hat{y}) &= \hat{y}(2\hat{y} - 1), \\ \hat{\phi}_4(\hat{x}, \hat{y}) &= 4\hat{x}\hat{y}, \\ \hat{\phi}_5(\hat{x}, \hat{y}) &= -4\hat{y}(\hat{x} + \hat{y} - 1), \\ \hat{\phi}_6(\hat{x}, \hat{y}) &= -4\hat{x}(\hat{x} + \hat{y} - 1),\end{aligned}\tag{21}$$

so that, according to (17), the elemental matrix for a generic triangle equals

$$A_{K_1}^{El} = \begin{bmatrix} 1 & \frac{1}{6} & \frac{1}{6} & 0 & -\frac{2}{3} & -\frac{2}{3} \\ \frac{1}{6} & \frac{1}{2} & 0 & 0 & 0 & -\frac{2}{3} \\ \frac{1}{6} & 0 & \frac{1}{2} & 0 & -\frac{2}{3} & 0 \\ 0 & 0 & 0 & \frac{8}{3} & -\frac{4}{3} & -\frac{4}{3} \\ -\frac{2}{3} & 0 & -\frac{2}{3} & -\frac{4}{3} & \frac{8}{3} & 0 \\ -\frac{2}{3} & -\frac{2}{3} & 0 & -\frac{4}{3} & 0 & \frac{8}{3} \end{bmatrix}\tag{22}$$

in the case of triangles of Type 1, or a suitable permutation in the case of triangles of Type 2. Despite the use of Lagrangian quadratic approximation, the stiffness matrix  $A_n = A_n(1, \Omega, \mathbb{P}_2)$  shows again a block tridiagonal structure  $A_n = \text{tridiag}(A_1, A_0, A_{-1})$  as in the linear case, the higher approximation stressing its influence just inside the blocks  $A_i$ . We might say that the quoted tridiagonal structure refers once again to triangles' vertices, while the internal structure is stressing the increased number of additional nodal points, that is, three in the case at hand. In fact, we observe in each block  $A_i$  a  $2 \times 2$  block structure as follows:

$$A_0 = \begin{bmatrix} B_0^{11} & B_0^{12} \\ (B_0^{12})^T & B_0^{22} \end{bmatrix}, \quad A_{-1} = \begin{bmatrix} 0 & 0 \\ B_{-1}^{21} & B_{-1}^{22} \end{bmatrix}, \quad A_1 = A_{-1}^T,\tag{23}$$

where the superscripts  $i, j$  in  $B_l^{ij}$  denote the position inside the  $2 \times 2$  block and the subscript  $l$  the belonging to the block  $A_l$ , so that

$$A_n = \begin{bmatrix} \begin{array}{cc|cc|cc} B_0^{11} & B_0^{12} & 0 & 0 & & \\ (B_0^{12})^T & B_0^{22} & B_{-1}^{21} & B_{-1}^{22} & & \\ \hline 0 & (B_{-1}^{21})^T & B_0^{11} & B_0^{12} & 0 & 0 \\ 0 & (B_{-1}^{22})^T & (B_0^{12})^T & B_0^{22} & B_{-1}^{21} & B_{-1}^{22} \\ \hline & & \ddots & & \ddots & \\ \hline & & 0 & (B_{-1}^{21})^T & B_0^{11} & B_0^{12} & 0 \\ & & 0 & (B_{-1}^{22})^T & (B_0^{12})^T & B_0^{22} & B_{-1}^{21} \\ \hline & & & & 0 & (B_{-1}^{21})^T & B_0^{11} \end{array} \\ \hline \end{bmatrix}.\tag{24}$$

More important, the very same structure depicted in (24), including the very same cutting in the lower right corner, appears in every block  $B_l^{ij}$  by considering suitable  $2 \times 2$  matrices as follows:

$$B_l^{ij} = \text{tridiag} \left( a_1^{B_l^{ij}}, a_0^{B_l^{ij}}, a_{-1}^{B_l^{ij}} \right), \quad l \in \{-1, 0, 1\}, \quad i, j \in \{1, 2\},$$

where

$$a_0^{B_0^{11}} = \begin{bmatrix} \frac{16}{3} & -\frac{4}{3} \\ -\frac{4}{3} & \frac{16}{3} \end{bmatrix}, \quad a_{-1}^{B_0^{11}} = \begin{bmatrix} 0 & 0 \\ -\frac{4}{3} & 0 \end{bmatrix}, \quad a_1^{B_0^{11}} = \left( a_{-1}^{B_0^{11}} \right)^T,$$

$$a_0^{B_0^{22}} = \begin{bmatrix} \frac{16}{3} & -\frac{4}{3} \\ -\frac{4}{3} & 4 \end{bmatrix}, \quad a_{-1}^{B_0^{22}} = \begin{bmatrix} 0 & 0 \\ -\frac{4}{3} & \frac{1}{3} \end{bmatrix}, \quad a_1^{B_0^{22}} = \left( a_{-1}^{B_0^{22}} \right)^T,$$

$$a_0^{B_0^{12}} = -\frac{4}{3}I_2, \quad a_{-1}^{B_0^{12}} = a_1^{B_0^{12}} = O_2,$$

$$a_0^{B_{-1}^{21}} = -\frac{4}{3}I_2, \quad a_{-1}^{B_{-1}^{21}} = a_1^{B_{-1}^{21}} = O_2,$$

$$a_0^{B_{-1}^{22}} = \begin{bmatrix} 0 & 0 \\ 0 & \frac{1}{3} \end{bmatrix}, \quad a_{-1}^{B_{-1}^{22}} = a_1^{B_{-1}^{22}} = O_2.$$

Thus, once again, just by taking into account that we are now facing a matrix-valued symbol, we can easily read the underlying symbol as follows:

$$f_{\mathbb{P}_2}(\theta_1, \theta_2) = f_{A_0}(\theta_1) + f_{A_{-1}}(\theta_1)e^{-i\theta_2} + f_{A_1}(\theta_1)e^{i\theta_2}, \quad (25)$$

with

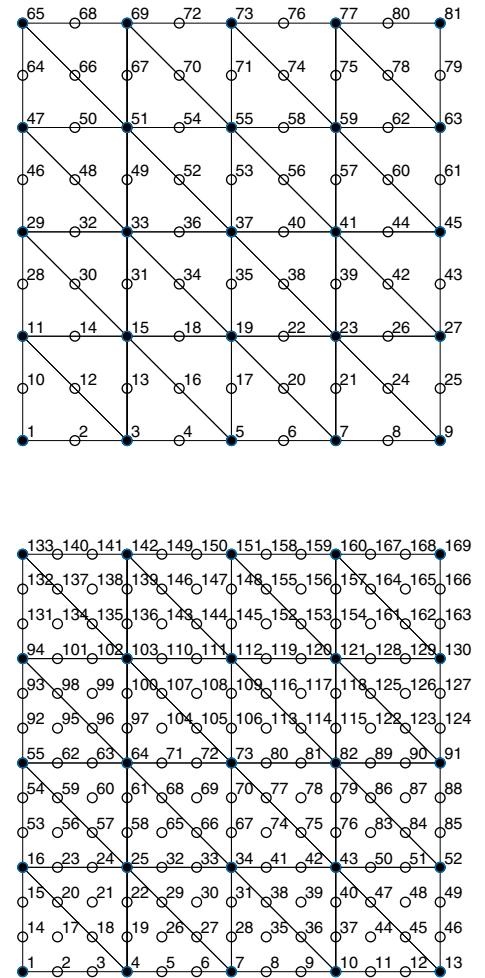
$$\begin{aligned} f_{A_0}(\theta_1) &= \begin{bmatrix} f_{B_0^{11}}(\theta_1) & f_{B_0^{12}}(\theta_1) \\ f_{(B_0^{12})^T}(\theta_1) & f_{B_0^{22}}(\theta_1) \end{bmatrix}, & f_{B_l^{ij}}(\theta_1) &= a_0^{B_l^{ij}} + a_{-1}^{B_l^{ij}}e^{-i\theta_1} + a_1^{B_l^{ij}}e^{i\theta_1}, \\ & & f_{(B_l^{ij})^T}(\theta_1) &= \overline{f_{B_l^{ij}}(\theta_1)}, \\ f_{A_{-1}}(\theta_1) &= \begin{bmatrix} 0 & 0 \\ f_{B_{-1}^{21}}(\theta_1) & f_{B_{-1}^{22}}(\theta_1) \end{bmatrix}, & f_{A_1}(\theta_1) &= \begin{bmatrix} 0 & f_{(B_{-1}^{21})^T}(\theta_1) \\ 0 & f_{(B_{-1}^{22})^T}(\theta_1) \end{bmatrix}. \end{aligned}$$

To sum up, we have a matrix-valued symbol  $f_{\mathbb{P}_2} : [-\pi, \pi]^2 \rightarrow \mathbb{C}^{4 \times 4}$  with

$$f_{\mathbb{P}_2}(\theta_1, \theta_2) = \left[ \begin{array}{cc|cc} \alpha & -\beta(1 + e^{i\theta_1}) & -\beta(1 + e^{i\theta_2}) & 0 \\ -\beta(1 + e^{-i\theta_1}) & \alpha & 0 & -\beta(1 + e^{i\theta_2}) \\ \hline -\beta(1 + e^{-i\theta_2}) & 0 & \alpha & -\beta(1 + e^{i\theta_1}) \\ 0 & -\beta(1 + e^{-i\theta_2}) & -\beta(1 + e^{-i\theta_1}) & \gamma + \frac{\beta}{2}(\cos(\theta_1) + \cos(\theta_2)) \end{array} \right] \quad (26)$$

with  $\alpha = 16/3$ ,  $\beta = 4/3$ , and  $\gamma = 4$ .

Finally, it is worth stressing that the stiffness matrix  $A_n(1, \Omega, \mathbb{P}_2)$  is a principal submatrix of a suitable permutation of the Toeplitz matrix  $T_n(f_{\mathbb{P}_2})$  defined according to (10). Indeed, the size of the two-level matrix  $A_n = A_n(1, \Omega, \mathbb{P}_2)$  is intrinsically odd both in inner and outer dimensions (see Theorem 2 and the explanation after Equation (13)), while  $T_n(f_{\mathbb{P}_2})$  has even corresponding dimensions: It is enough to cut every last row/column in each inner block, together with the last block with respect rows and columns, in order to obtain  $A_n$  from  $T_n(f_{\mathbb{P}_2})$ . In other words  $A_n$  is a special principal submatrix of  $T_n(f_{\mathbb{P}_2})$  according to the rule given in Theorem 2. As for the permutation we have just to consider the one defined by ordering nodal values as reported in Figure 3, where internal nodal values are grouped four by four. As a consequence the two matrix-sequences  $\{T_n(f_{\mathbb{P}_2})\}_n$  and  $\{A_n(1, \Omega, \mathbb{P}_2)\}_n$  share the same spectral distribution, that is, the same spectral symbol  $f_{\mathbb{P}_2}$ , by invoking Theorem 1. The following proposition holds.

**FIGURE 3** Nodal points reordering in  $\mathbb{P}_2$  and  $\mathbb{P}_3$  cases

**Proposition 2.** *The two matrix-sequences  $\{T_n(f_{\mathbb{P}_2})\}_n$  and  $\{A_n(1, \Omega, \mathbb{P}_2)\}_n$  are spectrally distributed as  $f_{\mathbb{P}_2}$  in the sense of Definition 1.*

As an immediate consequence of Proposition 2, we deduce a corollary regarding the clustering and localization of the spectra of  $\{T_n(f_{\mathbb{P}_2})\}_n$  and  $\{A_n(1, \Omega, \mathbb{P}_2)\}_n$ .

**Corollary 1.** *The range of  $f_{\mathbb{P}_2}$  is a weak cluster set for the spectra of the two matrix-sequences  $\{T_n(f_{\mathbb{P}_2})\}_n$  and  $\{A_n(1, \Omega, \mathbb{P}_2)\}_n$  in the sense of Definition 2. Furthermore, the convex hull of the range of  $f_{\mathbb{P}_2}$  contains all the eigenvalues of the involved matrices.*

*Proof.* The proof of the first part is a direct consequence of Proposition 2, taking into account Reference 36, Theorem 4.2 and observing that in this setting the standard range and the essential range coincide since  $f_{\mathbb{P}_2}$  is continuous (see also the end of Subsection 3.1). For the second part we observe that the result is known for Toeplitz matrices with Hermitian valued symbols:<sup>40</sup> Then the localization result for the eigenvalues of  $A_n(1, \Omega, \mathbb{P}_2)$  follows because  $A_n(1, \Omega, \mathbb{P}_2)$  is a principal submatrix of  $T_n(f_{\mathbb{P}_2})$  and since all the involved matrices are Hermitian. ■

### 5.3 | Case $k = 3$

In the case of cubic Lagrangian FE ( $k = 3$ ), by referring to the reference element (see Figure 2), we have the following shape functions

$$\hat{\varphi}_1(\hat{x}, \hat{y}) = -\frac{9}{2}\hat{x}^3 - \frac{27}{2}\hat{x}^2\hat{y} + 9\hat{x}^2 - \frac{27}{2}\hat{x}\hat{y}^2 + 18\hat{x}\hat{y} - \frac{11}{2}\hat{x} - \frac{9}{2}\hat{y}^3 + 9\hat{y}^2 - \frac{11}{2}\hat{y} + 1,$$

$$\begin{aligned}
\hat{\varphi}_2(\hat{x}, \hat{y}) &= \frac{\hat{x}}{2}(9\hat{x}^2 - 9\hat{x} + 2), \\
\hat{\varphi}_3(\hat{x}, \hat{y}) &= \frac{\hat{y}}{2}(9\hat{y}^2 - 9\hat{y} + 2), \\
\hat{\varphi}_4(\hat{x}, \hat{y}) &= \frac{9}{2}\hat{x}\hat{y}(3\hat{x} - 1), \\
\hat{\varphi}_5(\hat{x}, \hat{y}) &= \frac{9}{2}\hat{x}\hat{y}(3\hat{y} - 1), \\
\hat{\varphi}_6(\hat{x}, \hat{y}) &= -\frac{9}{2}\hat{y}(3\hat{y} - 1)(\hat{x} + \hat{y} - 1), \\
\hat{\varphi}_7(\hat{x}, \hat{y}) &= \frac{9}{2}\hat{y}(3\hat{x}^2 + 6\hat{x}\hat{y} - 5\hat{x} + 3\hat{y}^2 - 5\hat{y} + 2), \\
\hat{\varphi}_8(\hat{x}, \hat{y}) &= \frac{9}{2}\hat{x}(3\hat{x}^2 + 6\hat{x}\hat{y} - 5\hat{x} + 3\hat{y}^2 - 5\hat{y} + 2), \\
\hat{\varphi}_9(\hat{x}, \hat{y}) &= -\frac{9}{2}\hat{x}(3\hat{x} - 1)(\hat{x} + \hat{y} - 1), \\
\hat{\varphi}_{10}(\hat{x}, \hat{y}) &= -27\hat{x}\hat{y}(\hat{x} + \hat{y} - 1),
\end{aligned}$$

so that, according to (17), the elemental matrix for a generic triangle equals

$$A_{K_1}^{El} = \begin{bmatrix} \frac{17}{20} & -\frac{7}{80} & -\frac{7}{80} & -\frac{3}{40} & -\frac{3}{40} & \frac{3}{8} & -\frac{51}{80} & -\frac{51}{80} & \frac{3}{8} & 0 \\ -\frac{7}{80} & \frac{17}{40} & 0 & \frac{3}{80} & \frac{3}{80} & -\frac{3}{80} & -\frac{3}{80} & \frac{27}{80} & -\frac{27}{40} & 0 \\ -\frac{7}{80} & 0 & \frac{17}{40} & \frac{3}{80} & \frac{3}{80} & -\frac{27}{40} & \frac{27}{80} & -\frac{3}{80} & -\frac{3}{80} & 0 \\ -\frac{3}{40} & \frac{3}{80} & \frac{3}{80} & \frac{27}{80} & -\frac{27}{40} & \frac{27}{80} & \frac{27}{80} & \frac{27}{80} & -\frac{27}{16} & -\frac{81}{40} \\ -\frac{3}{40} & \frac{3}{80} & \frac{3}{80} & -\frac{27}{40} & \frac{27}{8} & -\frac{27}{16} & \frac{27}{80} & \frac{27}{80} & \frac{27}{80} & -\frac{81}{40} \\ \frac{3}{8} & -\frac{3}{80} & -\frac{27}{40} & \frac{27}{80} & -\frac{27}{16} & \frac{27}{8} & -\frac{27}{16} & 0 & 0 & 0 \\ -\frac{51}{80} & -\frac{3}{80} & \frac{27}{80} & \frac{27}{80} & \frac{27}{80} & -\frac{27}{16} & \frac{27}{8} & 0 & 0 & -\frac{81}{40} \\ -\frac{51}{80} & \frac{27}{80} & -\frac{3}{80} & \frac{27}{80} & \frac{27}{80} & 0 & 0 & \frac{27}{8} & -\frac{27}{16} & -\frac{81}{40} \\ \frac{3}{8} & -\frac{27}{40} & -\frac{3}{80} & -\frac{27}{16} & \frac{27}{8} & 0 & 0 & -\frac{27}{16} & \frac{27}{8} & 0 \\ 0 & 0 & 0 & -\frac{81}{40} & -\frac{81}{40} & 0 & -\frac{81}{40} & -\frac{81}{40} & 0 & \frac{81}{10} \end{bmatrix} \quad (27)$$

in the case of triangles of Type 1, or a suitable permutation in the case of triangles of Type 2. The stiffness matrix  $A_n = A_n(1, \Omega, \mathbb{P}_3)$  shows again a block tridiagonal structure  $A_n = \text{tridiag}(A_1, A_0, A_{-1})$  as in previous cases, the higher approximation stressing its influence just inside the blocks  $A_i$ . In fact, we observe in each block  $A_i$  a  $3 \times 3$  block structure as follows:

$$A_0 = \begin{bmatrix} B_0^{11} & B_0^{12} & B_0^{13} \\ (B_0^{12})^T & B_0^{22} & B_0^{23} \\ (B_0^{13})^T & B_0^{23} & B_0^{33} \end{bmatrix}, \quad A_{-1} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ B_{-1}^{31} & B_{-1}^{32} & B_{-1}^{33} \end{bmatrix}, \quad A_1 = A_{-1}^T. \quad (28)$$

More important, the very same structure appears in every block  $B_l^{ij}$  by considering suitable  $3 \times 3$  matrices and indeed we have

$$B_l^{ij} = \text{tridiag} \left( a_1^{B_l^{ij}}, a_0^{B_l^{ij}}, a_{-1}^{B_l^{ij}} \right), \quad l \in \{-1, 0, 1\}, \quad i, j \in \{1, 2, 3\},$$

where

$$a_0^{B_0^{11}} = \begin{bmatrix} \frac{81}{10} & -\frac{81}{40} & 0 \\ -\frac{81}{40} & \frac{27}{4} & -\frac{27}{16} \\ 0 & -\frac{27}{16} & \frac{27}{4} \end{bmatrix}, \quad a_{-1}^{B_0^{11}} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ -\frac{81}{40} & \frac{27}{80} & 0 \end{bmatrix}, \quad a_1^{B_0^{11}} = \left( a_{-1}^{B_0^{11}} \right)^T,$$



$$\begin{aligned}
a_0^{B_0^{22}} &= \begin{bmatrix} \frac{27}{4} & -\frac{81}{40} & \frac{27}{80} \\ -\frac{81}{40} & \frac{81}{10} & -\frac{81}{40} \\ \frac{27}{80} & -\frac{81}{40} & \frac{27}{4} \end{bmatrix}, \quad a_{-1}^{B_0^{22}} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ -\frac{27}{16} & 0 & 0 \end{bmatrix}, \quad a_1^{B_0^{22}} = \left(a_{-1}^{B_0^{22}}\right)^T, \\
a_0^{B_0^{33}} &= \begin{bmatrix} \frac{27}{4} & -\frac{27}{8} & \frac{57}{80} \\ -\frac{27}{8} & \frac{27}{4} & -\frac{21}{16} \\ \frac{57}{80} & -\frac{21}{16} & \frac{17}{5} \end{bmatrix}, \quad a_{-1}^{B_0^{33}} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ -\frac{21}{16} & \frac{57}{80} & -\frac{7}{40} \end{bmatrix}, \quad a_1^{B_0^{33}} = \left(a_{-1}^{B_0^{33}}\right)^T, \\
a_0^{B_0^{12}} &= \begin{bmatrix} -\frac{81}{40} & 0 & 0 \\ -\frac{27}{20} & -\frac{81}{40} & \frac{27}{80} \\ \frac{27}{80} & 0 & -\frac{27}{8} \end{bmatrix}, \quad a_{-1}^{B_0^{12}} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ \frac{27}{80} & 0 & 0 \end{bmatrix}, \quad a_1^{B_0^{12}} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & \frac{27}{80} \\ 0 & 0 & 0 \end{bmatrix}, \\
a_0^{B_0^{13}} &= \begin{bmatrix} 0 & 0 & 0 \\ \frac{27}{80} & \frac{27}{80} & -\frac{3}{40} \\ 0 & 0 & \frac{57}{80} \end{bmatrix}, \quad a_{-1}^{B_0^{13}} = O_3, \quad a_1^{B_0^{13}} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & \frac{3}{40} \\ 0 & 0 & -\frac{3}{80} \end{bmatrix}, \\
a_0^{B_0^{23}} &= \begin{bmatrix} -\frac{27}{16} & \frac{27}{80} & -\frac{3}{40} \\ 0 & -\frac{81}{40} & 0 \\ 0 & 0 & -\frac{21}{16} \end{bmatrix}, \quad a_{-1}^{B_0^{23}} = O_3, \quad a_1^{B_0^{23}} = \begin{bmatrix} 0 & 0 & \frac{3}{40} \\ 0 & 0 & 0 \\ 0 & 0 & -\frac{3}{80} \end{bmatrix}, \\
a_0^{B_{-1}^{31}} &= \begin{bmatrix} -\frac{81}{40} & \frac{27}{80} & 0 \\ 0 & -\frac{27}{16} & 0 \\ 0 & \frac{3}{40} & -\frac{21}{16} \end{bmatrix}, \quad a_{-1}^{B_{-1}^{31}} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & -\frac{3}{40} & 0 \end{bmatrix}, \quad a_1^{B_{-1}^{31}} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & -\frac{3}{80} \end{bmatrix}, \\
a_0^{B_{-1}^{32}} &= \begin{bmatrix} \frac{27}{80} & 0 & 0 \\ \frac{27}{80} & 0 & 0 \\ \frac{3}{40} & 0 & \frac{57}{80} \end{bmatrix}, \quad a_{-1}^{B_{-1}^{32}} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ -\frac{3}{40} & 0 & 0 \end{bmatrix}, \quad a_1^{B_{-1}^{32}} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & -\frac{3}{80} \end{bmatrix}, \\
a_0^{B_{-1}^{33}} &= \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ -\frac{3}{80} & -\frac{3}{80} & -\frac{7}{40} \end{bmatrix}, \quad a_{-1}^{B_{-1}^{33}} = O_3, \quad a_1^{B_{-1}^{33}} = \begin{bmatrix} 0 & 0 & -\frac{3}{80} \\ 0 & 0 & -\frac{3}{80} \\ 0 & 0 & 0 \end{bmatrix}.
\end{aligned}$$

Thus, once again, just by taking into account that we are now facing a matrix-valued symbol, we can easily read the underlying symbol as follows:

$$f_{\mathbb{P}_3}(\theta_1, \theta_2) = f_{A_0}(\theta_1) + f_{A_{-1}}(\theta_1)e^{-i\theta_2} + f_{A_1}(\theta_1)e^{i\theta_2}, \quad (29)$$

with

$$\begin{aligned}
f_{A_0}(\theta_1) &= \begin{bmatrix} f_{B_0^{11}}(\theta_1) & f_{B_0^{12}}(\theta_1) & f_{B_0^{13}}(\theta_1) \\ f_{(B_0^{12})^T}(\theta_1) & f_{B_0^{22}}(\theta_1) & f_{B_0^{23}}(\theta_1) \\ f_{(B_0^{13})^T}(\theta_1) & f_{(B_0^{23})^T}(\theta_1) & f_{B_0^{33}}(\theta_1) \end{bmatrix}, \quad f_{B_l^{ij}}(\theta_1) = a_0^{B_l^{ij}} + a_{-1}^{B_l^{ij}}e^{-i\theta_1} + a_1^{B_l^{ij}}e^{i\theta_1}, \\
&\quad f_{(B_l^{ij})^T}(\theta_1) = \overline{f_{B_l^{ij}}(\theta_1)}, \\
f_{A_{-1}}(\theta_1) &= \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ f_{B_{-1}^{31}}(\theta_1) & f_{B_{-1}^{32}}(\theta_1) & f_{B_{-1}^{33}}(\theta_1) \end{bmatrix}, \quad f_{A_1}(\theta_1) = \begin{bmatrix} 0 & 0 & f_{(B_{-1}^{31})^T}(\theta_1) \\ 0 & 0 & f_{(B_{-1}^{32})^T}(\theta_1) \\ 0 & 0 & f_{(B_{-1}^{33})^T}(\theta_1) \end{bmatrix}.
\end{aligned}$$

To sum up, we find the expression of  $f_{\mathbb{P}_3} : [-\pi, \pi]^2 \rightarrow \mathbb{C}^{9 \times 9}$  with

$$f_{\mathbb{P}_3}(\theta_1, \theta_2) = \begin{bmatrix} \alpha & -\frac{\alpha}{4} & -\frac{\alpha}{4}e^{i\theta_1} & -\frac{\alpha}{4} & 0 & 0 & -\frac{\alpha}{4}e^{i\theta_2} & 0 & 0 \\ -\frac{\alpha}{4} & \beta & -\frac{\beta}{4} + \frac{\beta}{20}e^{i\theta_1} & -\frac{\beta}{4} & -\frac{\alpha}{4} & \frac{\beta}{20}(1+e^{i\theta_1}) & \frac{\beta}{20}(1+e^{i\theta_2}) & \frac{\beta}{20} - \frac{\beta}{4}e^{i\theta_2} & f_{29} \\ -\frac{\alpha}{4}e^{-i\theta_1} & -\frac{\beta}{4} + \frac{\beta}{20}e^{-i\theta_1} & \beta & \frac{\beta}{20}(1+e^{-i\theta_1}) & 0 & -\frac{\beta}{2} & 0 & 0 & f_{39} \\ -\frac{\alpha}{4} & -\frac{\beta}{4} & \frac{\beta}{20}(1+e^{i\theta_1}) & \beta & -\frac{\alpha}{4} & \frac{\beta}{20} - \frac{\beta}{4}e^{i\theta_1} & -\frac{\beta}{4} + \frac{\beta}{20}e^{i\theta_2} & \frac{\beta}{20}(1+e^{i\theta_2}) & f_{49} \\ 0 & -\frac{\alpha}{4} & 0 & -\frac{\alpha}{4} & \alpha & -\frac{\alpha}{4} & 0 & -\frac{\alpha}{4} & 0 \\ 0 & \frac{\beta}{20}(1+e^{-i\theta_1}) & -\frac{\beta}{2} & \frac{\beta}{20} - \frac{\beta}{4}e^{-i\theta_1} & -\frac{\alpha}{4} & \beta & 0 & 0 & f_{69} \\ -\frac{\alpha}{4}e^{-i\theta_2} & \frac{\beta}{20}(1+e^{-i\theta_2}) & 0 & -\frac{\beta}{4} + \frac{\beta}{20}e^{-i\theta_2} & 0 & 0 & \beta & -\frac{\beta}{2} & f_{79} \\ 0 & \frac{\beta}{20} - \frac{\beta}{4}e^{-i\theta_2} & 0 & \frac{\beta}{20}(1+e^{-i\theta_2}) & -\frac{\alpha}{4} & 0 & -\frac{\beta}{2} & \beta & f_{89} \\ 0 & f_{29} & f_{39} & f_{49} & 0 & f_{69} & f_{79} & f_{89} & f_{99} \end{bmatrix} \quad (30)$$

where

$$\begin{aligned} f_{29} &= -\gamma(1 - e^{i\theta_1})(1 - e^{i\theta_2}), & f_{39} &= \delta - \frac{\gamma}{2}e^{i\theta_1} - \varepsilon e^{i\theta_2} - \frac{\gamma}{2}e^{-i\theta_1}e^{i\theta_2}, \\ f_{49} &= -\gamma(1 - e^{i\theta_1})(1 - e^{i\theta_2}), & f_{69} &= -\varepsilon - \frac{\gamma}{2}e^{i\theta_1} + \delta e^{i\theta_2} - \frac{\gamma}{2}e^{-i\theta_1}e^{i\theta_2}, \\ f_{79} &= \delta - \varepsilon e^{i\theta_1} - \frac{\gamma}{2}e^{i\theta_1}e^{-i\theta_2} - \frac{\gamma}{2}e^{i\theta_2}, & f_{89} &= -\varepsilon + \delta e^{i\theta_1} - \frac{\gamma}{2}e^{i\theta_1}e^{-i\theta_2} - \frac{\gamma}{2}e^{i\theta_2}, \\ f_{99} &= \zeta - 2\eta(\cos(\theta_1) + \cos(\theta_2)), \end{aligned}$$

and  $\alpha = 81/10$ ,  $\beta = 27/4$ ,  $\gamma = 3/40$ ,  $\delta = 57/80$ ,  $\varepsilon = 21/16$ ,  $\zeta = 17/5$ , and  $\eta = 7/40$ .

Finally, the stiffness matrix  $A_n = A_n(1, \Omega, \mathbb{P}_3)$  is a principal submatrix of a suitable permutation of the Toeplitz matrix  $T_n(f_{\mathbb{P}_3})$ . In order to obtain the stiffness matrix  $A_n$  from  $T_n(f_{\mathbb{P}_3})$ , it is enough to group internal nodal values nine by nine. Again by referring to Theorem 1, the following proposition holds.

**Proposition 3.** *The two matrix-sequences  $\{T_n(f_{\mathbb{P}_3})\}_n$  and  $\{A_n(1, \Omega, \mathbb{P}_3)\}_n$  are spectrally distributed as  $f_{\mathbb{P}_3}$  in the sense of Definition 1.*

As an immediate consequence of Proposition 3, we deduce a corollary regarding the clustering and localization of the spectra of  $\{T_n(f_{\mathbb{P}_3})\}_n$  and  $\{A_n(1, \Omega, \mathbb{P}_3)\}_n$  as well.

**Corollary 2.** *The range of  $f_{\mathbb{P}_3}$  is a weak cluster set for the spectra of the two matrix-sequences  $\{T_n(f_{\mathbb{P}_3})\}_n$  and  $\{A_n(1, \Omega, \mathbb{P}_3)\}_n$  in the sense of Definition 2. Furthermore, the convex hull of the range of  $f_{\mathbb{P}_2}$  contains all the eigenvalues of the involved matrices.*

*Proof.* The thesis follows with the same reasoning considered in Corollary 1. ■

## 6 | SYMBOL SPECTRAL ANALYSIS

We start the spectral analysis of symbols obtained in the previous section from a numerical point of view. As well known,  $f_{\mathbb{P}_1}$  shows a zero of order 2 in  $(0, 0)$ , while it is positive elsewhere. In the case  $k \geq 2$  the symbol is a matrix-valued function, so we consider an equispaced sampling in  $[-\pi, \pi]^2$  of the symbol and for each point we evaluate the  $k^2$  eigenvalues, ordering them in nondecreasing way. Thus,  $k^2$  surfaces  $s_i$ ,  $i = 1, \dots, k^2$ , are defined, and the  $i$ th eigenvalue in a given point of the sampling being the value of the surface  $s_i$  in such a point. In Table 1, the minimal and maximal values of each surface  $s_i$ ,  $i = 1, \dots, k^2$ , are reported (for a comparison among the surfaces obtained by using the eigenvalues of the considered matrix-sequence and the corresponding surfaces obtained by properly sampling the symbol of the same matrix-sequence see the subsequent Figures 4–7).

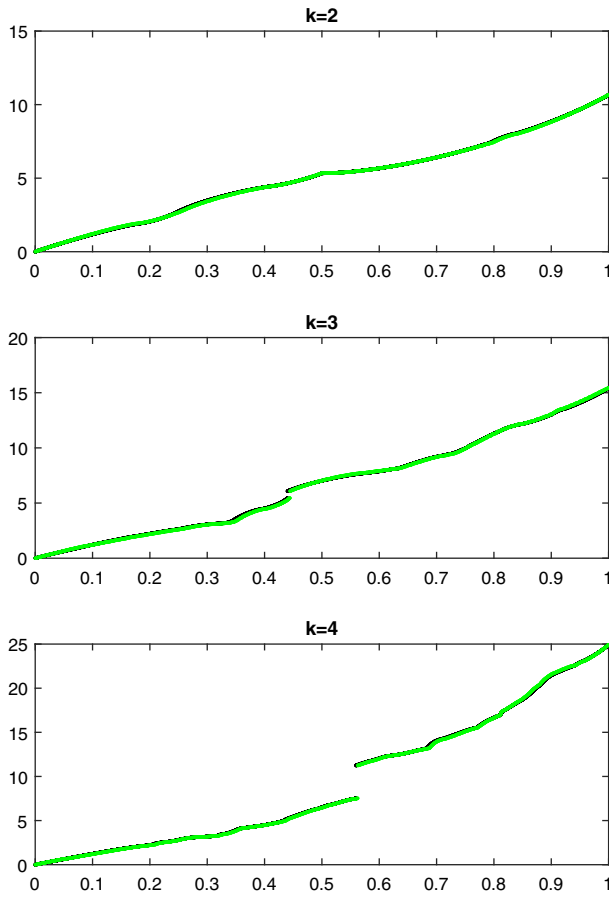
In the case  $k = 2$ , it is worth stressing that the chosen sorting of the eigenvalues influences the surfaces definition, the minimal value of the  $i$ th surface being lower of the maximal value of the  $(i - 1)$ th surface: This implies that the union of the ranges of the eigenvalue functions of the symbol produces a connected set which is a cluster for the spectra of the given matrix-sequence. When  $k = 3$  the union of the ranges of the first four surfaces is well separated from the union of the remaining five surfaces and hence the cluster is divided into two subclusters in the sense of Definition 2. In the case  $k = 4$ , the union of the ranges of the first nine surfaces is well separated from the union of the remaining seven surfaces and consequently, as in the case of  $k = 3$ , the cluster is divided into two subclusters.

However, there is a phenomenon which is expected and it is independent of the value of  $k$ : only the first surface reaches zero as minimum, while all the other surfaces are strictly positive everywhere.

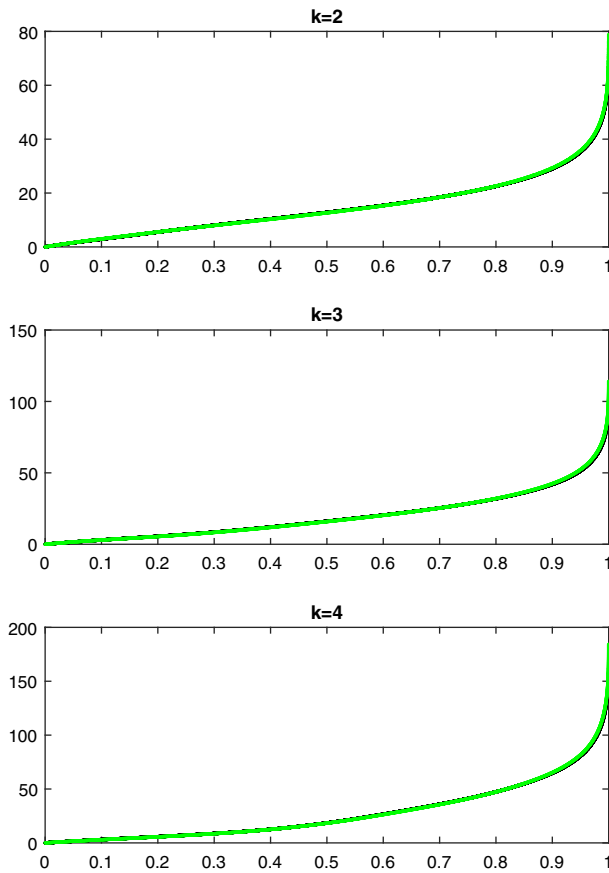
**TABLE 1** Minimum and maximum of surfaces  $s_i$ ,  $i = 1, \dots, k^2$ 

$i$	$\min(s_i)$	$\operatorname{argmin}(s_i)$	$\max(s_i)$	$\operatorname{argmax}(s_i)$
$k = 1$				
1	0	(0, 0)	8	$(-\pi, -\pi)$
$k = 2$				
1	-2.122988181725368e-17	(0, 0)	2.666666666666667e+00	$(-\pi, -\pi)$
2	2.666666666666667e+00	(0, $-\pi$ )	5.333333333333330e+00	(0, 0)
3	5.333333333333325e+00	$(-u, -u)$	7.415403750411773e+00	(0, $-\pi$ )
4	5.33333333333333e+00	$(-\pi, -\pi)$	1.066666666666667e+01	(0, 0)
$k = 3$				
1	-2.947870832408496e-16	(0, 0)	1.752299219210445e+00	$(-\pi, -\pi)$
2	1.077001420967619e+00	$(-\pi, 0)$	2.649326100400095e+00	(0, 0)
3	2.02499999999999e+00	$(\pi, -\pi)$	3.37499999999998e+00	(0, 0)
4	2.417725227846304e+00	$(-\pi, -\pi)$	5.473900873539699e+00	(0, 0)
5	6.07499999999998e+00	(0, 0)	8.150826984062711e+00	$(-v, v)$
6	6.07500000000001e+00	(0, 0)	9.45000000000005e+00	$(-\pi, -\pi)$
7	8.10000000000000e+00	$(-\pi, 0)$	1.145177302606021e+01	(0, 0)
8	1.01250000000000e+01	$(-\pi, -\pi)$	1.306461248424784e+01	$(-\pi, 0)$
9	1.21500000000001e+01	(0, 0)	1.542003087979332e+01	$(-\pi, -\pi)$
$k = 4$				
1	8.665811124242140e-15	(0, 0)	1.154132889535501e+00	$(\pi, -\pi)$
2	6.091179158637314e-01	$(-\pi, 0)$	2.028216383055356e+00	$(-\pi, -\pi)$
3	1.183562035003593e+00	$(w, -w)$	2.278229389751864e+00	$(-\pi, 0)$
4	1.228189268889777e+00	$(-\pi, -\pi)$	2.796565325232735e+00	$(z, \pi/10)$
5	2.706589845271391e+00	(0, 0)	3.280192294561743e+00	$(-\pi, -\pi)$
6	3.100532625333826e+00	(0, 0)	4.876190476190478e+00	$(\pi, -\pi)$
7	4.086126343234464e+00	$(a, b)$	5.001164336911962e+00	$(c, -c)$
8	4.923102258884252e+00	$(d, -d)$	6.507154754218933e+00	(0, $-\pi$ )
9	6.351060427971650e+00	$(e, e)$	7.524544180802064e+00	$(f, -f)$
10	1.124369260315089e+01	$(-\pi, 0)$	1.277464393947364e+01	$(\pi, -\pi)$
11	1.221231815611003e+01	$(-g, g)$	1.319265448651837e+01	(0, $-\pi$ )
12	1.312292935024724e+01	$(h, h)$	1.551857649581396e+01	$(i, -i)$
13	1.403908928314477e+01	$(\pi, -\pi)$	1.715087064330462e+01	$(l, -l)$
14	1.715307867863427e+01	$(l, -l)$	2.041132693707404e+01	$(-\pi, -\pi)$
15	1.987073604165514e+01	(0, 0)	2.261907577519149e+01	(0, $\pi$ )
16	2.236934933910699e+01	$(m, -m)$	2.492211941947813e+01	(0, 0)

Note:  $u = -7.351326809400116e - 01$ ,  $v = 2.500707752257475e + 00$ ,  
 $w = 3.053628059289279e + 00$ ,  $z = 1.734159144781565e + 00$ ,  $a = 2.576105975943630e - 01$ ,  
 $b = -2.224247598741574e + 00$ ,  $c = 2.896548426609789e + 00$ ,  $d = 1.507964473723100e + 00$ ,  
 $e = 1.043008760991811e + 00$ ,  $f = 2.161415745669778e + 00$ ,  $g = 1.627344994559513e + 00$ ,  
 $h = 2.796017461694915e + 00$ ,  $i = 7.099999397112930e - 01$ ,  $l = 9.550441666912972e - 01$ ,  
 $m = 2.519557308179014e + 00$ .

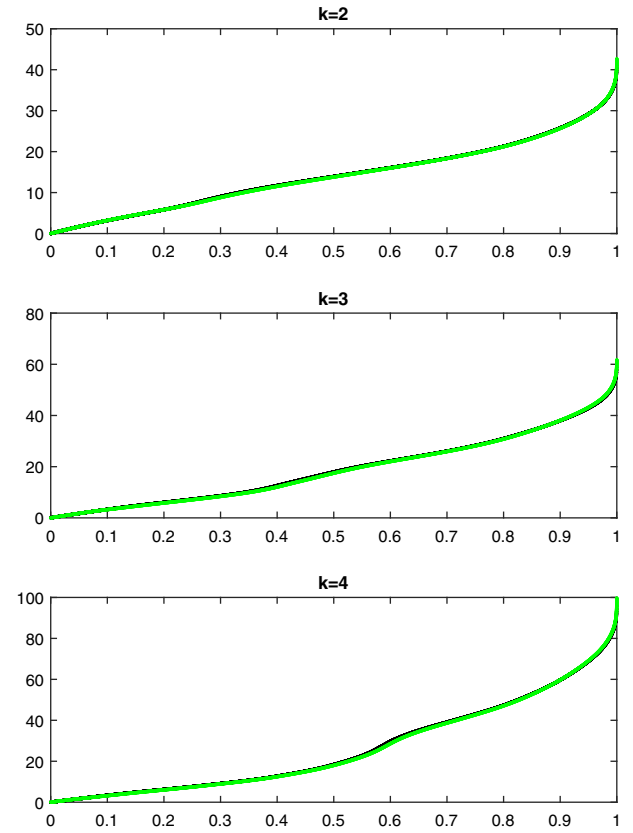


**FIGURE 4** Ordered equispaced samplings of  $\lambda_j(a(x,y)f_{\mathbb{P}_k}(\theta))$ ,  $j = 1, \dots, k^2$  (green dots) and ordered eigenvalues  $\lambda_l(A_n(a, \Omega, \mathbb{P}_k))$  with  $a(x,y) \equiv 1$

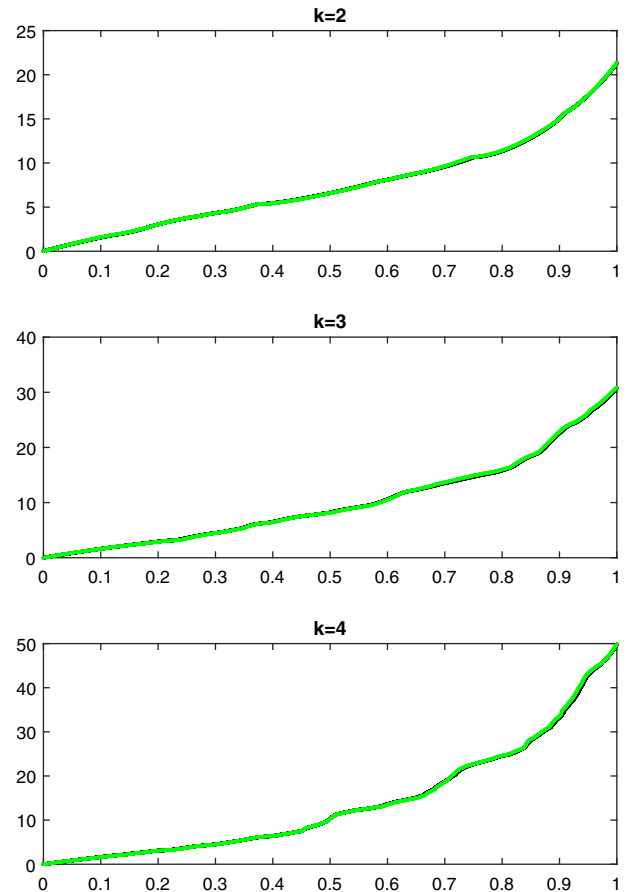


**FIGURE 5** Ordered equispaced samplings of  $\lambda_j(a(x,y)f_{\mathbb{P}_k}(\theta))$ ,  $j = 1, \dots, k^2$  (green dots) and ordered eigenvalues  $\lambda_l(A_n(a, \Omega, \mathbb{P}_k))$  with  $a(x,y) = e^{x+y}$ .

**FIGURE 6** Ordered equispaced samplings of  $\lambda_j \left( a(x, y) f_{\mathbb{P}_k}(\theta) \right)$ ,  $j = 1, \dots, k^2$  (green dots) and ordered eigenvalues  $\lambda_l(A_n(a, \Omega, \mathbb{P}_k))$  with  $a(x, y) = 1 + 2\sqrt{x} + y$



**FIGURE 7** Ordered equispaced samplings of  $\lambda_j \left( a(x, y) f_{\mathbb{P}_k}(\theta) \right)$ ,  $j = 1, \dots, k^2$  (green dots) and ordered eigenvalues  $\lambda_l(A_n(a, \Omega, \mathbb{P}_k))$  with  $a(x, y) = 1$  if  $y \geq x$  and  $a(x, y) = 2$  otherwise



Now we give a general result regarding the main features of the involved symbols, with the proof in various cases, including both  $\mathbb{P}_k$  and  $\mathbb{Q}_k$  Finite Element approximations.

**Theorem 4.** *Given the symbols  $f_{\mathbb{P}_k}, f_{\mathbb{Q}_k}$  in dimension  $d \geq 1$ , the following statements hold true. For every  $f \in \{f_{\mathbb{P}_k}, f_{\mathbb{Q}_k}\}$ , setting*

$$\lambda_1(f(\boldsymbol{\theta})) \leq \cdots \leq \lambda_{k^d}(f(\boldsymbol{\theta})),$$

we obtain

1.  $f(\mathbf{0})\mathbf{e} = 0$ ,  $\mathbf{e}$  vector of all ones,  $k \geq 1$ ;
2. there exist constants  $C_1, C_2 > 0$  (dependent on  $f$ ) such that

$$C_1 \sum_{j=1}^d (2 - 2 \cos(\theta_j)) \leq \lambda_1(f(\boldsymbol{\theta})) \leq C_2 \sum_{j=1}^d (2 - 2 \cos(\theta_j)); \quad (31)$$

3. there exist constants  $m, M > 0$  (dependent on  $f$ ) such that

$$0 < m \leq \lambda_j(f(\boldsymbol{\theta})) \leq M, \quad j = 2, \dots, k^d. \quad (32)$$

For  $f_{\mathbb{Q}_k}$ , the proof is given for every  $k, d \geq 1$  (for  $d = 1$  we notice again that  $f_{\mathbb{P}_k} \equiv f_{\mathbb{Q}_k}$ ). For  $f_{\mathbb{P}_k}$ , the proof is given for  $d = 2$  and  $k = 2, 3$ .

*Remark 1.* Although in the case of  $\mathbb{Q}_k$  Finite Elements the analysis of the symbol  $f_{\mathbb{Q}_k}$  given in Reference 13 and in Theorem 4 is general, for the  $\mathbb{P}_k$  Finite Elements in dimension  $d > 1$  there is still room for a substantial improvement of the analysis and this will be a target in future researches.

*Proof.* Case  $\mathbb{Q}_k$  Finite Elements: any  $k \geq 1, d = 1$ .

Claims 2. and 3. have been proved in Theorem 8 and Corollary 1 in Reference 13. Here, we prove Claim 1.: As first thing we recall that the relation  $f(\mathbf{0})\mathbf{e} = 0$ , with  $\mathbf{e}$  vector of all ones and  $k \geq 1$ , is equivalent to say that every row of  $f(\mathbf{0})$  is a vector having rowsum equal to zero. We now show the latter feature. Taking into consideration the notations in Section 4, we have

$$(f(\mathbf{0})\mathbf{e})_s = \sum_{j=1}^k (f(\mathbf{0}))_{sj} = \sum_{j=1}^k (K_0 + K_1 + K_1^T)_{sj}, \quad s = 1, \dots, k.$$

We first observe that the Lagrange polynomial interpolating the constant 1 is exactly equal to 1, by the uniqueness of the interpolant. Therefore  $\sum_{j=0}^k L_j = 1$ ,  $\left(\sum_{j=0}^k L_j\right)' = \sum_{j=0}^k L'_j = 0$ , and hence, for  $1 \leq s \leq k-1$ ,

$$\begin{aligned} \sum_{j=1}^k (f(\mathbf{0}))_{sj} &= \sum_{j=1}^k \langle L'_j, L'_s \rangle + \langle L'_0, L'_s \rangle \\ &= \left\langle \sum_{j=0}^k L'_j, L'_s \right\rangle \\ &= \langle 0, L'_s \rangle = 0. \end{aligned}$$

Finally, for  $s = k$  we have

$$\sum_{j=1}^k (f(\mathbf{0}))_{kj} = \sum_{j=1}^k \langle L'_j, L'_k \rangle + \langle L'_0, L'_k \rangle + \langle L'_0, L'_k \rangle + \sum_{j=1}^k \langle L'_0, L'_j \rangle$$



$$\begin{aligned}
&= \sum_{j=0}^k \langle L'_j, L'_k \rangle + \sum_{j=0}^k \langle L'_0, L'_j \rangle \\
&= \langle \sum_{j=0}^k L'_j, L'_k \rangle + \langle L'_0, \sum_{j=0}^k L'_j \rangle = 0,
\end{aligned}$$

and consequently we conclude that  $f(\mathbf{0})\mathbf{e} = 0$ .

Case  $\mathbb{Q}_k$  Finite Elements: any  $k \geq 1$ , any  $d \geq 2$ . Claim 1. is a direct consequence of the proof for  $\mathbb{Q}_k$  Finite Elements,  $k \geq 1$ , and  $d = 1$ , given its tensorial structure (see formula (5.1) in Reference 13): In reality, it is sufficient to observe that  $x \otimes y$  has rowsum equal to zero if and only if either  $x$  or  $y$  has rowsum equal to zero, with any  $x, y$  complex vectors of any size. The case of more than two vectors can be handled by an inductive argument. Furthermore, Claims 2. and 3. are contained in Section 5.1 in Reference 13.

Case  $\mathbb{P}_k$  Finite Elements:  $k = 2$ ,  $d = 2$ . Claim 1. follows by direct check of the zero rowsum property from the expression of the symbol  $f_{\mathbb{P}_k}$  in (26), taking into account  $\theta = (0, 0)$  and the numerical values of the involved parameters.

Now, since the determinant of a matrix is the product of its eigenvalues and since  $f$  is bounded in infinity norm, in order to prove Claim 2. and 3. with  $d = 2$  it is sufficient to show that:

- A.  $\det(f(\theta)) \sim \sum_{j=1}^2 (2 - 2 \cos(\theta_j))$ ,  $\theta = (\theta_1, \theta_2)$ ,
  - B. there exists  $C > 0$  such that  $\lambda_2(f(\theta)) \geq C > 0$ ,
- with  $f = f_{\mathbb{P}_2}$ .

We remind that the relation A. means there exist  $C_1, C_2 > 0$  such that

$$C_1 \sum_{j=1}^2 (2 - 2 \cos(\theta_j)) \leq \det(f_{\mathbb{P}_2}(\theta)) \leq C_2 \sum_{j=1}^2 (2 - 2 \cos(\theta_j))$$

uniformly in the domain  $(\theta_1, \theta_2) \in [-\pi, \pi]^2$ .

By direct computation, we find

$$\begin{aligned}
\det(f_{\mathbb{P}_2}(\theta)) &= C' (-2 \cos(\theta_1) - 2 \cos(\theta_2) - \cos(\theta_1) \cos(\theta_2) + 5) \\
&= C' \left( \sum_{j=1}^2 (2 - 2 \cos(\theta_j)) + (1 - \cos(\theta_1) \cos(\theta_2)) \right) \\
&\geq C' \left( \sum_{j=1}^2 (2 - 2 \cos(\theta_j)) \right)
\end{aligned}$$

with  $C' = 4096/81$ , being  $-1 \leq \cos(\theta_1) \cos(\theta_2) \leq 1$  for all  $(\theta_1, \theta_2) \in [-\pi, \pi]^2$ . Thus,  $C_1 = C'$ .

Furthermore, for  $c = 1/2$  it holds  $1 - \cos(\theta_1) \cos(\theta_2) \leq c \left( \sum_{j=1}^2 (2 - 2 \cos(\theta_j)) \right)$  for all  $(\theta_1, \theta_2) \in [-\pi, \pi]^2$ . Thus,  $C_2 = 3C'/2$ .

Finally, let  $\lambda_1 \leq \lambda_2 \leq \lambda_3 \leq \lambda_4$  be the eigenvalues of the Hermitian matrix  $f_{\mathbb{P}_2}(\theta)$  and let  $\mu_1 \leq \mu_2 \leq \mu_3$  be the eigenvalues of the principal submatrix  $g(\theta)$  chosen as  $g(\theta) = (f_{\mathbb{P}_2}(\theta))_{i,j=2}^4$ .

Since the approximation matrices of problem (1) are all positive definite due to

- coerciveness of the continuous problem,
- the use of Galerkin techniques such as the Finite Elements,

it follows that the symbol

$$f_{\mathbb{P}_2}(\theta)$$

of the related matrix-sequence has to be Hermitian nonnegative definite which means that  $\lambda_1 \geq 0$  on the whole definition domain. By contradiction if  $\lambda_1$  is negative in a set of positive measure then, by the distribution results,<sup>8,9</sup> many eigenvalues of the approximation matrices would be negative for a matrix size large enough and this is impossible.

By using the interlacing theorem, we have

$$\lambda_1 \leq \mu_1 \leq \lambda_2 \leq \mu_2 \leq \lambda_3 \leq \mu_3 \leq \lambda_4, \quad (33)$$

with  $\lambda_1$  equal to zero at  $\theta = \mathbf{0}$  and positive elsewhere. By direct computation of the determinant we find that  $\det(g(\theta)) > 0$  and hence, taking into account that  $g(\theta)$  is continuous and positive definite on the compact square  $[-\pi, \pi]^2$ , we conclude that all the eigenvalues of  $g(\theta)$  are strictly positive and continuous on  $[-\pi, \pi]^2$  that is  $\mu_j > 0$  for  $j = 1, 2, 3$ . Thus, using (34), we conclude  $\lambda_2 \geq \mu_1 \geq \min_{\theta \in [-\pi, \pi]^2} \mu_1 > 0$ .

Case  $\mathbb{P}_k$  Finite Elements:  $k = 3, d = 2$ . As in the case  $k = 2$ , Claim 1. follows by direct inspection from the expression of the symbol  $f_{\mathbb{P}_k}$  in (31), taking into account  $\theta = (0, 0)$  and the numerical values of the involved parameters.

In order to prove Claims 2. and 3. we follow the very same steps as for the case  $k = 2$ , that is, we prove **A.** and **B.** with  $f = f_{\mathbb{P}_3}$ .

By direct computation we have

$$\begin{aligned} \det(f_{\mathbb{P}_3}(\theta)) &= a(-\cos(\theta_2)\cos^2(\theta_1) - \cos(\theta_1)\cos^2(\theta_2) + 4\cos^2(\theta_1) + 4\cos^2(\theta_2) \\ &\quad - 80\cos(\theta_1)\cos(\theta_2) - 195\cos(\theta_1) - 195\cos(\theta_2) + 464), \end{aligned}$$

where  $a = 205,891,132,094,649/81,920,000,000$ . We write  $\det(f_{\mathbb{P}_3}(\theta))$  in the form

$$\det(f_{\mathbb{P}_3}(\theta)) = a \left( h(\theta) + \frac{195}{2} \sum_{j=1}^2 (2 - 2\cos(\theta_j)) \right),$$

where

$$h(\theta) = -\cos(\theta_2)\cos^2(\theta_1) - \cos(\theta_1)\cos^2(\theta_2) + 4\cos^2(\theta_1) + 4\cos^2(\theta_2) - 80\cos(\theta_1)\cos(\theta_2) + 74.$$

Since  $-\cos^2(\theta_k) \leq -\cos(\theta_j)\cos^2(\theta_k)$  and  $1 - \cos(\theta_1)\cos(\theta_2) \geq 0$ , we obtain

$$\begin{aligned} h(\theta) &\geq 3\cos^2(\theta_1) + 3\cos^2(\theta_2) - 80\cos(\theta_1)\cos(\theta_2) + 74 \\ &\geq 3(\cos(\theta_1) - \cos(\theta_2))^2 - 74\cos(\theta_1)\cos(\theta_2) + 74 \\ &\geq 0, \end{aligned}$$

which implies directly  $\det(f_{\mathbb{P}_3}(\theta)) \geq C_1 \sum_{j=1}^2 (2 - 2\cos(\theta_j))$  with  $C_1 = \frac{195}{2}a$ .

On the other side, taking into account  $\cos^2(\theta_j) \leq 1, j = 1, 2$ , we deduce

$$\begin{aligned} h(\theta) &= \cos^2(\theta_1)(4 - \cos(\theta_2)) + \cos^2(\theta_2)(4 - \cos(\theta_1)) - 80\cos(\theta_1)\cos(\theta_2) + 74 \\ &\leq 8 - \cos(\theta_1) - \cos(\theta_2) - 80\cos(\theta_1)\cos(\theta_2) + 74 \\ &\leq 2 - \cos(\theta_1) - \cos(\theta_2) + 80(1 - \cos(\theta_1)\cos(\theta_2)) \\ &\leq \frac{81}{2} \sum_{j=1}^2 (2 - 2\cos(\theta_j)). \end{aligned}$$

Due to the relation  $1 - \cos(\theta_1)\cos(\theta_2) \leq \frac{1}{2}(4 - 2\cos(\theta_1) - 2\cos(\theta_2))$ , as already observed in the case  $k = 2$ , we find  $\det(f_{\mathbb{P}_3}(\theta)) \leq 138a \sum_{j=1}^2 (2 - 2\cos(\theta_j))$  with  $C_2 = 138a$ .

Finally, let  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_9$  be the eigenvalues of the Hermitian matrix  $f_{\mathbb{P}_3}(\theta)$ , and let  $g(\theta) = (f_{\mathbb{P}_3}(\theta))_{i,j=1}^8$  be the principal submatrix and  $\mu_1 \leq \dots \leq \mu_8$  its eigenvalues.

Since the approximation matrices of problem (1) are all positive definite due to

- coerciveness of the continuous problem,
- the use of Galerkin techniques such as the Finite Elements,

it follows that the symbol  $f_{\mathbb{P}_3}(\theta)$  of the related matrix-sequence has to be Hermitian nonnegative definite which means that  $\lambda_1 \geq 0$  on the whole definition domain. By contradiction if  $\lambda_1$  is negative in a set of positive measure then, by the distribution results,<sup>8,9</sup> many eigenvalues of the approximation matrices would be negative for a matrix size large enough and this is impossible.

By using the interlacing theorem, we have

$$\lambda_1 \leq \mu_1 \leq \lambda_2 \leq \mu_2 \leq \dots \leq \mu_8 \leq \lambda_9, \quad (34)$$

with  $\lambda_1$  equal to zero at  $\theta = \mathbf{0}$  and positive elsewhere. By direct computation we find  $\det(g(\theta)) = \prod_{j=1}^8 \mu_j > 0$  so that, taking into account that  $g(\theta)$  is continuous and positive definite on the compact square  $[-\pi, \pi]^2$ , we deduce  $\mu_j > 0$  for all  $j = 1, \dots, 8$ . Consequently, by (35), we conclude  $\lambda_2 \geq \mu_1 \geq \min_{\theta \in [-\pi, \pi]^2} \mu_1 > 0$ . ■

## 6.1 | Extremal eigenvalues and conditioning

As already observed, direct consequences of Proposition 2, Corollary 1, Proposition 3, and Corollary 2 are that the sequences of Finite Element matrices are distributed as the symbol  $f$  and that the union of the ranges of the eigenvalue functions of  $f$  represent a cluster for their spectra, while the convex hull of the union of the ranges of the eigenvalue functions of  $f$  contains all the eigenvalues of the involved matrices.

On the other hand, Theorem 4 gives information on the analytical properties of  $f$ , which are relevant for giving results on the extreme eigenvalues and the asymptotic conditioning.

Indeed, from Theorem 4, we know that the minimal eigenvalue function of  $f$  behaves as the symbol of the standard Finite Difference Laplacian, while the other eigenvalue functions are well separated from zero and bounded. Furthermore, thanks to the analysis in Reference 40, the fact that the minimal eigenvalue of  $f$  has a zero of order two implies that

- the minimal eigenvalue goes to zero as  $N^{-2/d}$ ,
- the maximal eigenvalue converges from below to the maximum of the maximal eigenvalue function of  $f$
- and hence the conditioning of the involved matrices grow asymptotically exactly as  $N^{2/d}$ ,

with  $N$  being the global matrix size (see also the argument in Section 5.1 in Reference 13 and Reference 41).

## 7 | THE CASE OF VARIABLE COEFFICIENTS AND NON-CARTESIAN DOMAINS

When the diffusion coefficient  $a(\mathbf{x})$  in (1) is not constant, the structure of the stiffness matrix is no longer Toeplitz, but somehow the Toeplitz character is hidden in an asymptotic sense and indeed the sequence of matrices  $\{A_n(a, \Omega, \mathbb{P}_k)\}_n$  approximating (1) can be spectrally treated with the help of the GLT technology with  $k = 1, 2, 3$ .

Below we report the essentials of the steps for computing the spectral symbol.

**Step 1.** If  $\Omega = (0, 1)^d$ ,  $d \geq 1$ , then  $\{A_n(a, \Omega, \mathbb{P}_k)\}_n$  can be written as a sequence of principal submatrices of a linear combination of products involving the multilevel block Toeplitz sequence generated by  $f_{\mathbb{P}_k}$ , the diagonal sampling sequence of  $a(\mathbf{x})I_{N(k,d)}$ , and zero distributed sequences. The use of items **GLT 1.–GLT 3.**, combined with Theorem 1, leads to the conclusion

$$\{A_n(a, \Omega, \mathbb{P}_k)\}_n \sim_{\sigma, \lambda} a(\mathbf{x})f_{\mathbb{P}_k}(\theta), \quad \mathbf{x} \in (0, 1)^d, \quad \theta \in [-\pi, \pi]^d. \quad (35)$$

**Step 2.** If  $\Omega$  is Peano–Jordan measurable, then without loss of generality, we assume  $\Omega \subset \Omega_d = (0, 1)^d$  and  $d \geq 2$ . Hence  $\{A_n(a, \Omega, \mathbb{P}_k)\}_n$  can be seen, up to zero-distributed sequences, as a sequence of principal submatrices of  $\{A_n(\hat{a}, \Omega, \mathbb{P}_k)\}_n$ , where  $\hat{a}$  is equal to  $a$  on the domain  $\Omega$  and it is identically zero in the complement  $\Omega_d \setminus \Omega$ . In this way we are reduced to Step 1. and the use of a reduction argument (see Section 6 in Reference 6 and Section 3.1.4 in Reference 7) implies the distribution result

$$\{A_n(a, \Omega, \mathbb{P}_k)\}_{n \sim \sigma, \lambda} a(\mathbf{x}) f_{\mathbb{P}_k}(\theta), \quad \mathbf{x} \in \Omega, \quad \theta \in [-\pi, \pi]^d. \quad (36)$$

The rest of the section is now devoted to show that the predictions in (36) and (37) are numerically confirmed. Indeed, in the constant coefficient case, we plotted the surface of the different eigenvalue functions  $\lambda_j(f_{\mathbb{P}_k}(\theta))$ ,  $j = 1, \dots, k^2$ ,  $k = 1, 2, 3, 4$ , and this was technically possible because the functions are all bivariate as  $\theta \in [-\pi, \pi]^2$ .

In the variable coefficient case, the visualization is substantially more involved, since the symbol is  $a(\mathbf{x}) f_{\mathbb{P}_k}(\theta)$  and hence the eigenvalue functions  $\lambda_j(a(\mathbf{x}) f_{\mathbb{P}_k}(\theta))$ ,  $j = 1, \dots, k^2$ ,  $k = 1, 2, 3, 4$ , are all functions in four variables as  $\mathbf{x} \in \Omega$ ,  $\theta \in [-\pi, \pi]^2$ . Consequently, for visualization purposes, we choose a different technique: for a fixed  $k$  and for a fixed matrix size, we make an ordering (nondecreasing) of all the eigenvalues of  $A_n(a, \Omega, \mathbb{P}_k)$  and we take the same ordering (nondecreasing) of the values given by an equispaced sampling of all the functions  $\lambda_j(a(\mathbf{x}) f_{\mathbb{P}_k}(\theta))$ ,  $j = 1, \dots, k^2$ .

As it can be seen from Figures 4 to 7, all concerning the case  $\Omega = (0, 1)^2$ , the match is perfect showing that the distribution result in (36) is fully confirmed with  $a(x, y) = 1$ ,  $a(x, y) = e^{x+y}$ ,  $a(x, y) = 1 + 2\sqrt{x} + y$ ,  $a(x, y) = 1$  if  $y \geq x$  and  $a(x, y) = 2$  otherwise.

We have four relevant remarks.

- As a general observation, the graph of the ordered equispaced sampling of  $\lambda_j(a(\mathbf{x}) f_{\mathbb{P}_k}(\theta))$ ,  $j = 1, \dots, k^2$ , represent a monotone rearrangement<sup>42</sup> of the different eigenvalue functions and this global rearrangement is a fortiori a univariate monotone function.
- When  $a(\mathbf{x}) \equiv 1$ , the symbol in (36) reduces to  $f_{\mathbb{P}_k}$  and we observe jumps which correspond to the existence of an index  $l$  such that

$$\max \lambda_l(f_{\mathbb{P}_k}(\theta)) < \lambda_{l+1}(f_{\mathbb{P}_k}(\theta)),$$

with  $1 \leq l \leq k^2 - 1$ ,  $k = 1, 2, 3, 4$ . In other words, the rearranged function has a few discontinuity points. When  $a(\mathbf{x})$  is not constant such a phenomenon disappears, since the range of all eigenvalue functions becomes wider and all the ranges intersect (there is not a range not intersecting at least another range). Furthermore, beside the latter smoothing effect due to  $a(\mathbf{x})$ , it is worthwhile observing that the regularity of the diffusion coefficient does not affect the qualitative behavior of the eigenvalue distribution. In fact, the reconstruction given by the symbol of the eigenvalues of  $A_n(a, \Omega, \mathbb{P}_k)$  is accurate in all the considered examples and, more specifically, the figures look very similar independently of the fact that the diffusion coefficient is smooth ( $a(\mathbf{x}) = e^{x+y}$  in Figure 5), or is  $C^0$  but not  $C^1$  ( $a(\mathbf{x}) = 1 + 2\sqrt{x} + y$  in Figure 6), or is discontinuous ( $a(\mathbf{x}) = 1$  if  $y \geq x$  and  $a(\mathbf{x}) = 2$  otherwise, in Figure 7).

- In all Figures 5–7, the matrix size is quite moderate (of the order of  $10^4$ ), showing that the spectral distribution effects, which represent an asymptotic property, can be already visualized for small orders of the considered matrices.
- We did not show any figure regarding the distribution formula (37) just because the check has been done and there is no difference with respect to the case of  $\Omega = (0, 1)^2$  as in (36).

## 8 | PRECONDITIONING AND COMPLEXITY ISSUES

Finally, we consider a few numerical experiments on preconditioning. First of all, the interest in solving the constant coefficient case  $a(x, y) \equiv 1$  refers to its use in optimally preconditioning the nonconstant coefficient case whenever  $a$  is smooth enough (see References 43 and 44 for the case  $k = 1$ ). This is evident from Table 2, where we report the number of iterations required by PCG applied to  $A_n(a, \Omega, \mathbb{P}_k)$ , with  $a(x, y) = \exp(x + y)$  in the case of preconditioning with  $P_n(a) = \tilde{D}_n^{1/2}(a) A_n(1, \Omega, \mathbb{P}_k) \tilde{D}_n^{1/2}(a)$  with  $\tilde{D}_n(a) = D_n(a) D_n^{-1}(1)$ , where  $D_n(a)$  is the main diagonal of  $A_n(a, \Omega, \mathbb{P}_k)$  and  $D_n(1)$  is the main diagonal of  $A_n(1, \Omega, \mathbb{P}_k)$ .

Therefore, we now focus our attention on the case  $a(x, y) \equiv 1$ . Taking into account the Toeplitz nature of the matrices at hand, our aim is to preliminarily test a classical preconditioner as the Circulant one, clearly by considering the Strang correction to deal with its singularity. More precisely we will consider as preconditioner the Circulant matrix generated by the very same function  $f_{\mathbb{P}_k}(\theta)$  plus the correction  $h^2 ee^T$ ,  $e$  being the vector of all ones and  $h$  the constant triangle edge. In

**TABLE 2** Number of PCG's iterations to reach convergence with respect to relative residual less than  $1.e-6$ , Preconditioner  $P_n(a)$ ,  $a(x, y) = \exp(x + y)$ 

$k = 2$		$k = 3$		$k = 4$	
$N$	$P_n(a)$	$N$	$P_n(a)$	$N$	$P_n(a)$
49	4	121	5	225	5
225	3	529	4	961	5
961	3	2,209	4	3,969	4
3,969	3	9,025	4	16,129	4
16,129	3	36,481	4	65,025	4

**TABLE 3** Number of PCG's iterations to reach convergence with respect to relative residual less than  $1.e-6$ , case  $k = 2$ 

$k = 2$					
$N$	$P = I$		$P = IC$		$P = C_S$
64	26	19	11	10	14
256	47	42	18	17	19
1,024	90	86	32	31	26
4,096	174	170	56	55	38
16,384	336	331	98	96	53

**TABLE 4** Number of outliers  $n_{\text{out}}$  (eigenvalues not belonging to  $(1 - \varepsilon, 1 + \varepsilon)$  with  $\varepsilon = 1.e - 1$ ) and their percentage with respect the dimension

$k = 2$					
$N$	$n_{\text{out}}$	%	$n_{\text{out}}$	%	
64	27	$4.2e-1$	27	$4.2e-1$	
256	59	$2.3e-1$	55	$2.1e-1$	
1,024	123	$1.2e-1$	111	$1.1e-1$	
4,096	251	$6.1e-2$	225	$5.5e-2$	

*Note:* The second and third columns refer to the Toeplitz case, the fourth and fifth columns to the case of the FEM matrix.

the even columns of Table 3 (case  $k = 2$ ), we report the number of PCG iterations required to solve the system with Toeplitz matrix  $T_n(f_{\mathbb{P}_k})$ , in the case of no preconditioning ( $P_n = I_n$ ), preconditioning by the incomplete Cholesky factorization, and by the Circulant  $C_n(f_{\mathbb{P}_k})$  plus the Strang correction, respectively. To this end, it is worth stressing that we have to consider the dimension of the Toeplitz/Circulant matrix fitting with its natural dimension with respect to the symbol size. Therefore, when instead we want to solve the system with the FEM matrix  $A_n(1, \Omega, \mathbb{P}_k)$ , principal submatrix of the matrix  $T_n(f_{\mathbb{P}_k})$ , we need to match its dimension with the one previously considered for the Circulant preconditioner. We obtain that goal just by imposing boundary conditions to  $T_n(f_{\mathbb{P}_k})$ , but keeping the size unchanged. The related numerical results are reported in odd columns of Table 3. In both cases, the number of required iterations increases as the dimension increases. No significant difference in even or odd column is observed in the case of no preconditioning or incomplete Cholesky preconditioning (the results are slightly better in the second case). In the case of the Circulant preconditioning, we observe a clear worsening in the effectiveness when the preconditioner is applied not to the Toeplitz matrix, but to its principal submatrix plus boundary conditions, though the iteration growth rate seems smaller than the one observed in the case of incomplete Cholesky preconditioning. Furthermore, as expected from the theory, a weak cluster around 1 is observed (see Table 4).

In Tables 5 and 6, the same numerical experiments are reported in the case  $k = 3$  and  $k = 4$ . The numerical behavior seems to be substantially of the same type, independently of the parameter  $k$ , also in reference to the weak cluster phenomenon observed for  $k = 2$  in Table 4.

**$k = 3$** 

$N$	$P = I$		$P = IC$		$P = C_S$	
144	45	39	16	15	19	30
576	82	78	28	27	26	42
2,304	159	155	50	48	36	61
9,216	306	301	86	84	49	90

**TABLE 5** Number of PCG's iterations to reach convergence with respect to relative residual less than  $1.e-6$ , case  $k = 3$  **$k = 4$** 

$N$	$P = I$		$P = IC$		$P = C_S$	
256	74	66	10	10	23	38
1,024	134	129	18	17	31	57
4,096	261	254	31	30	43	81
16,384	502	490	54	51	61	116

**TABLE 6** Number of PCG's iterations to reach convergence with respect to relative residual less than  $1.e-6$ , case  $k = 4$ 

## 9 | CONCLUDING REMARKS

We considered a class of elliptic partial differential equations with Dirichlet boundary conditions, where the operator is  $\text{div}(-a(\mathbf{x})\nabla \cdot)$  with  $a$  continuous and positive on  $\overline{\Omega}$ . For the numerical approximation we have chosen the classical  $\mathbb{P}_k$  Finite Element method, in the case of Friedrichs–Keller triangulations, leading to sequence of matrices of increasing size. The new results concern the spectral analysis of the resulting matrix-sequences both in the direction of the global distribution in the Weyl sense and of the asymptotic conditioning. We considered in detail the case of constant coefficients and we have given a brief account in the more involved case of variable coefficients. The mathematical tools stem out from the Toeplitz technology and from the rather new theory of GLT matrix-sequences. Numerical results are shown for a practical evidence of the theoretical findings.

Several open problems remain and here we make a short list of the most relevant ones:

- it would be valuable to find a unified formula for the symbols of the  $\mathbb{P}_k$  over a  $d$  dimensional cube for every  $k, d$  as done for the case of  $\mathbb{Q}_k$  discretizations;<sup>13</sup>
- the analysis in the variable coefficient case has to be completed, but this seems now a reasonable target given the new findings in the theory of GLT matrix-sequences (see References 9 and 45 and references therein);
- as expected in any multilevel setting,<sup>46</sup> the standard preconditioning procedures fail to be optimal: It would be of great value an extension to this context of the multigrid techniques developed for Toeplitz structures with scalar-valued symbols. Of course the proof of optimality is the final goal in this numerical setting.

## ACKNOWLEDGEMENTS

The authors warmly thank the anonymous referees for their careful reading and for the several critical suggestions, which have greatly improved the presentation and the content of the article. The work of S.Serra-Capizzano and C.Tablino-Possio is supported by GNCS. This work does not have any conflicts of interest.

## ORCID

Cristina Tablino-Possio  <https://orcid.org/0000-0003-1424-2767>

## REFERENCES

- Quarteroni A. Numerical models for differential problems. Milan, Italia: Springer-Verlag, 2014.
- Ciarlet PG. The finite element method for Elliptic problems. Philadelphia, PA: SIAM Publisher, 2002.
- Schwab C. p- and hp-finite element methods: Theory and applications in solid and fluid mechanics. Oxford, UK: Clarendon Press, 1998.



4. Brezzi F, Fortin M. Mixed and hybrid finite element methods. New York, NY: Springer-Verlag, 1991.
5. Böttcher A, Silbermann B. Introduction to Large Truncated Toeplitz Matrices. New York, NY: Springer-Verlag, 1999.
6. Serra-Capizzano S. Generalized locally Toeplitz sequences: Spectral analysis and applications to discretized partial differential equations. *Linear Algebra Appl.* 2003;366:371–402.
7. Serra-Capizzano S. The GLT class as a generalized Fourier analysis and applications. *Linear Algebra Appl.* 2006;419-1:180–233.
8. Garoni C, Serra-Capizzano S. Generalized locally Toeplitz sequences: Theory and applications. Vol I. Cham: Springer, 2017.
9. Garoni C, Serra-Capizzano S. Generalized locally Toeplitz sequences: Theory and applications. Vol II. Cham: Springer, 2018.
10. Beckermann B, Serra-Capizzano S. On the asymptotic spectrum of finite element matrix sequences. *SIAM J Numer Anal.* 2007;45-2:746–769.
11. Morozov S, Serra-Capizzano S, Tyrtysnikov E. How to extend the application scope of GLT-sequences. TR Dept Inf Tech, Division Scientific Comput. 2018;13:1–30. <http://www.it.uu.se/research/publications/reports/2018-013/>.
12. Dorostkar A, Neytcheva M, Serra-Capizzano S. Spectral analysis of coupled PDEs and of their Schur complements via generalized locally Toeplitz sequences in 2D. *Comput Methods Appl Mech Eng.* 2016;309:74–105.
13. Garoni C, Serra-Capizzano S, Sesana D. Spectral analysis and spectral symbol of d-variate  $\mathbf{Q}_p$  Lagrangian FEM stiffness matrices. *SIAM J Matrix Anal Appl.* 2015;36-3:1100–1128.
14. Tilli P. A note on the spectral distribution of Toeplitz matrices. *Linear Multilin Algebra.* 1998;45:147–159.
15. Saad Y. Iterative methods for sparse linear systems. 2nd ed. Philadelphia, PA: SIAM, 2003.
16. Aricò A, Donatelli M, Serra-Capizzano S. V-cycle optimal convergence for certain (multilevel) structured linear systems. *SIAM J Matrix Anal Appl.* 2004;26:186–214.
17. Beckermann B, Kuijlaars ABJ. Superlinear convergence of conjugate gradients. *SIAM J Numer Anal.* 2001;39:300–329.
18. Donatelli M, Garoni C, Manni C, Serra-Capizzano S, Speleers H. Robust and optimal multi-iterative techniques for IgA Galerkin linear systems. *Comput Methods Appl Mech Eng.* 2015;284:230–264.
19. Donatelli M, Garoni C, Manni C, Serra-Capizzano S, Speleers H. Robust and optimal multi-iterative techniques for IgA collocation linear systems. *Comput Methods Appl Mech Eng.* 2015;284:1120–1146.
20. Donatelli M, Garoni C, Manni C, Serra-Capizzano S, Speleers H. Symbol-based multigrid methods for Galerkin B-spline isogeometric analysis. *SIAM J Numer Anal.* 2017;55-1:31–62.
21. Fiorentino G, Serra S. Multigrid methods for Toeplitz matrices. *Calcolo.* 1991;28:283–305.
22. Mazza M, Ratnani A, Serra Capizzano S. Spectral analysis and spectral symbol for the 2D curl-curl (stabilized) operator with applications to the related iterative solutions. *Math Comput.* 2019;88-317:1155–1188.
23. Cottrell JA, Hughes TJR, Bazilevs Y. Isogeometric analysis: Toward integration of CAD and FEA. Chichester, NH: John Wiley & Sons, 2009.
24. Garoni C, Manni C, Pelosi F, Serra-Capizzano S, Speleers H. On the spectrum of stiffness matrices arising from isogeometric analysis. *Numer Math.* 2014;127:751–799.
25. Donatelli M, Garoni C, Manni C, Serra-Capizzano S, Speleers H. Spectral analysis and spectral symbol of matrices in isogeometric collocation methods. *Math Comput.* 2016;85-300:1639–1680.
26. Garoni C, Speleers H, Ekström S-E, Reali A, Serra-Capizzano S, Hughes TJR. Symbol-based analysis of finite element and isogeometric B-spline discretizations of eigenvalue problems: Exposition and review. *Arch Comput Meth Eng.* 2019;26(5):1639–1690. <https://doi.org/10.1007/s11831-018-9295-y>.
27. Engl HW, Hanke M, Neubauer A. Regularization of inverse problems. Dordrecht, Netherlands: Kluwer Academic Publishers, 2000.
28. Donatelli M, Molteni M, Pennati V, Serra-Capizzano S. Multigrid methods for cubic spline solution of two point (and 2D) boundary value problems. *Appl Numer Math.* 2016;104:15–29.
29. Benedusi P, Garoni C, Krause R, Li X, Serra-Capizzano S. Space-time FE-DG discretization of the anisotropic diffusion equation in any dimension: The spectral symbol. *SIAM J Matrix Anal Appl.* 2018;39-3:1383–1420.
30. Cottrell JA, Reali A, Bazilevs Y, Hughes TJR. Isogeometric analysis of structural vibrations. *Comput Methods Appl Mech Eng.* 2006;195:5257–5296.
31. Hughes TJR, Evans JA, Reali A. Finite element and NURBS approximations of eigenvalue, boundary-value, and initial-value problems. *Comput Methods Appl Mech Eng.* 2014;272:290–320.
32. Hughes TJR, Reali A, Sangalli G. Duality and unified analysis of discrete approximations in structural dynamics and wave propagation: Comparison of p-method finite elements with k-method NURBS. *Comput Methods Appl Mech Eng.* 2008;197:4104–4124.
33. Reali A. An isogeometric analysis approach for the study of structural vibrations. *J Earthq Eng.* 2006;10:1–30.
34. Ferrari P, Rahla R, Tablino-Possio C, Belhaj S, Serra-Capizzano S. Multigrid for  $\mathbf{Q}_k$  Finite Element Matrices using a (block) Toeplitz symbol approach. *Mathematics.* 2020;8-1:5. doi:10.3390/math8010005; <https://www.mdpi.com/2227-7390/8/1/5>.
35. Bhatia R. Matrix analysis. New York, NY: Springer-Verlag, 1997.
36. Golinskii L, Serra-Capizzano S. The asymptotic properties of the spectrum of nonsymmetrically perturbed Jacobi matrix sequences. *J Approx Th.* 2007;144-1:84–102.
37. Tyrtyshnikov E. A unifying approach to some old and new theorems on distribution and clustering. *Linear Algebra Appl.* 1996;232:1–43.
38. Garoni C, Serra-Capizzano S, Sesana D. Block locally Toeplitz sequences: Construction and properties. Paper presented at: Proceedings of the Workshop Structured Matrices Analysis, Algorithms and Applications - Cortona (AR - Italy) September 4-8, (2019):25-58; Springer INDAM Series 30.

39. Garoni C, Serra-Capizzano S, Sesana D. Block generalized locally toeplitz sequences: Topological construction, spectral distribution results, and star-algebra structure. Paper presented at: Proceedings of the Workshop Structured Matrices Analysis, Algorithms and Applications - Cortona (AR - Italy) September 4–8, 2017, Springer INDAM Series 30; (2019); 59–79.
40. Serra-Capizzano S. Asymptotic results on the spectra of block Toeplitz preconditioned matrices. *SIAM J Matrix Anal Appl.* 1998;20-1:31–44.
41. Serra-Capizzano S. On the extreme eigenvalues of Hermitian (block) Toeplitz matrices. *Linear Algebra Appl.* 1998;270:109–129.
42. Kawohl B. Rearrangements and convexity of level sets in PDE, lecture notes in mathematics. Berlin, Germany: Springer-Verlag, 1985;p. 1150.
43. Serra-Capizzano S, TablinoPossio C. Finite element matrix sequences: The case of rectangular domains. *Numer Alg.* 2001;28:309–327.
44. Russo A, Serra-Capizzano S, Tablino-Possio C. Quasi-optimal preconditioners for finite element approximations of diffusion dominated convection diffusion equations on (nearly) equilateral triangle meshes. *Numer Linear Algebra Appl.* 2015;22-1:123–144.
45. Garoni C, Mazza M, Serra-Capizzano S. Block generalized locally Toeplitz sequences: From the theory to the applications. *Axioms.* 2018;7:49.
46. Serra-Capizzano S, Tyrtysnikov E. Any circulant-like preconditioner for multilevel matrices is not superlinear. *SIAM J Matrix Anal Appl.* 1999;21-2:431–439.

## SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of this article.

**How to cite this article:** Rahla RI, Serra-Capizzano S, Tablino-Possio C. Spectral analysis of  $\mathbb{P}_k$  Finite Element matrices in the case of Friedrichs–Keller triangulations via Generalized Locally Toeplitz technology. *Numer Linear Algebra Appl.* 2020;e2302. <https://doi.org/10.1002/nla.2302>