

## ON THE ASYMPTOTIC SPECTRUM OF FINITE ELEMENT MATRIX SEQUENCES\*

BERNHARD BECKERMANN<sup>†</sup> AND STEFANO SERRA-CAPIZZANO<sup>‡</sup>

**Abstract.** We derive a new formula for the asymptotic eigenvalue distribution of stiffness matrices obtained by applying  $P_1$  finite elements with standard mesh refinement to the semielliptic PDE of second order in divergence form  $-\nabla(K\nabla^T u) = f$  on  $\Omega$ ,  $u = g$  on  $\partial\Omega$ . Here  $\Omega \subset \mathbb{R}^2$ , and  $K$  is supposed to be piecewise continuous and pointwise symmetric semipositive definite. The symbol describing this asymptotic eigenvalue distribution depends on the PDE, but also both on the numerical scheme for approaching the underlying bilinear form and on the geometry of triangulation of the domain. Our work is motivated by recent results on the superlinear convergence behavior of the conjugate gradient method, which requires the knowledge of such asymptotic eigenvalue distributions for sequences of matrices depending on a discretization parameter  $h$  when  $h \rightarrow 0$ . We compare our findings with similar results for the finite difference method which were published in recent years. In particular we observe that our sequence of stiffness matrices is part of the class of generalized locally Toeplitz sequences for which many theoretical tools are available. This enables us to derive some results on the conditioning and preconditioning of such stiffness matrices.

**Key words.** finite element methods, matrix sequence, asymptotic eigenvalue distribution

**AMS subject classifications.** 65F10, 65N22, 15A18, 15A12, 47B65

**DOI.** 10.1137/05063533X

**1. Introduction and statement of the main results.** Consider the semielliptic PDE of second order in divergence form

$$(1) \quad -\nabla(K\nabla^T u) = f \quad \text{on } \Omega, \quad u = g \quad \text{on } \partial\Omega,$$

where  $\Omega \subset \mathbb{R}^2$  is a bounded open “smooth” set (say, with piecewise  $\mathcal{C}^1$  boundary), and  $K : \Omega \mapsto \mathbb{R}^{2 \times 2}$  is piecewise continuous in  $\Omega$  and symmetric semipositive definite at each point of  $\Omega$ . In this paper we are interested in describing the asymptotic distribution of eigenvalues of the matrix of coefficients obtained by approximating the above elliptic PDE by  $P_1$  finite elements in the case where the position of the vertices can be described by some mapping.

The task of finding the asymptotic eigenvalue distribution is motivated by some recent results on the (superlinear) convergence behavior for the method of conjugate gradients (CG) [4, 5, 6]: a discretization of (1) for some sequence of stepsizes  $h$  tending to zero leads to a sequence of systems of linear equations  $A_n x_n = b_n$  with  $A_n$  some symmetric positive definite matrix of order  $n$ , where of course  $n$  depends on  $h$  and tends to  $\infty$  for  $h \rightarrow 0$ . The CG method is a popular method for solving such systems, and its convergence properties have been analyzed by many authors (see, e.g., [3, 41]). For instance, one has a simple upper bound for the CG error in the energy norm in

---

\*Received by the editors July 6, 2005; accepted for publication (in revised form) November 1, 2006; published electronically April 13, 2007.

<http://www.siam.org/journals/sinum/45-2/63533.html>

<sup>†</sup>Laboratoire de Mathématiques Paul Painlevé, UMR CNRS 8524, Université des Sciences et Technologies de Lille, F-59655 Villeneuve d’Ascq, France (Bernhard.Beckermann@univ-lille1.fr). This author’s work was supported in part by INTAS network NeCCA 03-51-6637, and in part by the grant FAR 2004 of the University of Como.

<sup>‡</sup>Dipartimento di Fisica e Matematica, Università dell’Insubria, Via Valleggio 11, 22100 Como, Italy (stefano.serrac@uninsubria.it, serra@mail.dm.unipi.it). This author’s work was supported in part by MIUR grant 2002014121 and in part by the Department of Mathematics, Université des Sciences et Technologies de Lille.

terms of the spectral condition number of  $A_n$ , that is, the ratio of the largest divided by the smallest eigenvalue of  $A_n$ ; see, e.g., [24, (6.106)]. Both for finite difference and finite element approximations, asymptotics for the smallest eigenvalue of  $A_n$  in terms of  $h$  and the smallest eigenvalue of the differential operator of (1) are known; see, for instance, [20]. By elementary means one also obtains upper bounds for the largest eigenvalue, and hence upper bounds for the CG error.

However, the (linear) upper bound based on the condition number is usually quite rough, especially in the range of superlinear convergence of CG. This superlinear convergence behavior is observed numerically to be quite pronounced in the context of discretized elliptic problems in  $\geq 2$  dimensions, in particular for small stepsizes  $h$ . Here CG convergence is known to be governed by the distribution of the spectrum  $\Lambda(A_n)$  of  $A_n$ , which at least for very simple model problems can be computed explicitly. Roughly speaking, superlinear CG convergence occurs if the eigenvalue distribution of  $A_n$  is far from being a worst case eigenvalue distribution. This qualitative rule of thumb has been known already for some time, but has been quantified only recently in [4, 5, 6]: here the authors give asymptotic error estimates for CG in terms of the asymptotic eigenvalue distribution of  $(A_n)_{n \geq 0}$ , namely the so-called *asymptotic spectrum* defined as follows.

A sequence of matrices  $(A_n)_{n \geq 0}$ ,  $A_n$  Hermitian of order  $n$  with spectrum  $\Lambda(A_n) \subset \mathbb{R}$ , is said to have an *asymptotic spectrum* given by some measure  $\sigma$  if for all functions  $f \in \mathcal{C}_c(\mathbb{R})$  (i.e., continuous with compact support) there holds

$$(2) \quad \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{\lambda \in \Lambda(A_n)} f(\lambda) = \int f(\lambda) d\sigma(\lambda),$$

where each eigenvalue is counted according to its multiplicity (and hence  $\sigma$  is a probability measure supported on the extended real line  $\mathbb{R} = \mathbb{R} \cup \{\pm\infty\}$ ). In the case where the limit (2) exists and takes the form

$$(3) \quad \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{\lambda \in \Lambda(A_n)} f(\lambda) = \int_D f(\omega(t)) \frac{dt}{m(D)}$$

with a domain  $D \subset \mathbb{R}^d$  having finite Lebesgue measure  $m(D) > 0$ , the function  $\omega$  will be referred to as the *symbol* of  $(A_n)$ .

The probably most classical example of sequences of matrices having an asymptotic spectrum is given by Hermitian Toeplitz matrices  $A_n = (t_{j-k})_{j,k=1,\dots,n}$  obtained from the Fourier coefficients of the Lebesgue integrable generating function  $\omega(s) = \sum_{j \in \mathbb{Z}} t_j e^{ijs}$ ,  $i^2 = -1$ ; see, for instance, [8] and references therein. Here the symbol coincides with the generating function, and  $D = (-\pi, \pi)$ .

In the present paper, the matrices  $A_n$  will result from the same approximation process when using different (decreasing) stepsizes, and thus one might expect that the sequence of matrices  $(A_n)$  has an asymptotic spectrum. Indeed, in case of finite difference discretization for differential operators, explicit formulas for an asymptotic spectrum have been given in [23, 38, 33, 26] (one-dimensional setting) and [31, 32, 30, 28, 35] (two-dimensional and multidimensional setting). Each time, the underlying symbol includes information on the coefficients and the domain of the PDE and information on the discretization schemes for the derivatives. To our knowledge, results for finite element approximations are still lacking (except for some preliminary results in [26, 31]).

Before stating our results on stiffness matrices for finite elements in subsection 1.2, we first recall in subsection 1.1 some known examples of asymptotic spectra in the finite difference case.

**1.1. The case of finite difference discretizations.** Consider the discretization of the one-dimensional boundary value problem

$$\begin{cases} -\frac{d}{dx} \left( k(x) \frac{d}{dx} u(x) \right) = f(x), & x \in (0, 1), \\ u(0), u(1) \text{ given numbers,} \end{cases}$$

on a uniformly spaced grid using centered finite differences of precision order 2 and minimal bandwidth. The resulting linear systems are of tridiagonal type with coefficient matrices  $(A_n)$  having entries which are weighted samples of  $k$ :

$$(4) \quad A_n = \begin{bmatrix} k_{\frac{1}{2}} + k_{\frac{3}{2}} & -k_{\frac{3}{2}} & & & & & \\ -k_{\frac{3}{2}} & k_{\frac{3}{2}} + k_{\frac{5}{2}} & -k_{\frac{5}{2}} & & & & \\ & -k_{\frac{5}{2}} & \ddots & \ddots & & & \\ & & \ddots & \ddots & & & \\ & & & & -k_{\frac{2n-1}{2}} & & \\ & & & -k_{\frac{2n-1}{2}} & k_{\frac{2n-1}{2}} + k_{\frac{2n+1}{2}} & & \end{bmatrix},$$

with  $k_t = k(t \cdot h)$ ,  $h = (n + 1)^{-1}$ . When  $k(x) \equiv 1$ , the matrix  $A_n$  reduces to the Toeplitz matrix  $T_n(a)$  of size  $n$ ,

$$(5) \quad T_n(a) = \begin{bmatrix} 2 & -1 & & & & \\ -1 & 2 & -1 & & & \\ & -1 & \ddots & \ddots & & \\ & & \ddots & \ddots & -1 & \\ & & & -1 & 2 & \end{bmatrix},$$

generated by  $a(s) = 2 - 2 \cos(s)$ : note that the numbers  $-1, 2, -1$  are the (nonzero) Fourier coefficients  $c_1, c_0, c_{-1}$  of  $a$  and represent also the stencil of the finite difference formula. This latter statement is not a coincidence: if we change the stencil (for instance, in order to obtain more precise discretization schemes), then we obtain Toeplitz matrices generated by a new function  $a$  having Fourier coefficients given by the entries of this new stencil [33]. A well-known fact from the theory of Toeplitz matrices is that  $(T_n(a))_n$  has an asymptotic spectrum given by  $\omega(s) = a(s)$  with  $D = [-\pi, \pi]$ ; see, for instance, the seminal work by Grenander and Szegö [17]. In the more general case of variable coefficients, it follows from the locally Toeplitz analysis of [38] that the matrices  $A_n$  of (4) have an asymptotic spectrum given by the symbol

$$\omega(x, s) = k(x)a(s)$$

with  $D = (0, 1) \times [-\pi, \pi]$  (see also [23]). We observe that the result is in some sense natural since the samplings of  $k$  move along the diagonals of  $A_n$  smoothly (if  $k$  is smooth), and therefore also the algebraic structure of  $A_n$  looks like a Toeplitz if we restrict our attention to a local portion of the matrix: this nice algebraic behavior has a natural counterpart in the global spectral behavior. As in the constant coefficient case, the change of the discretization scheme, i.e., of the stencil, will change only

the function  $a$  in the symbol (compare [33] and [38]). Finally, we observe that the matrices  $(A_n)$  are essentially of the same type as those which one encounters when dealing with sequences of orthogonal polynomials with varying coefficients. Here again locally Toeplitz tools have been used for finding the distribution of the zeros of the considered orthogonal polynomials under very weak assumptions (only measurability) on the regularity of the coefficients [22] (see also [40]).

A further variation which could be considered in the discretization of the above one-dimensional boundary value problem is the use of nonequispaced grids. Indeed, if the new grid of size  $n$  is obtained as the image under a map  $\phi : [0, 1] \mapsto [0, 1]$  of a uniform grid of the same size  $n$  or if the new grid can be approximated in this way (see, e.g., [35, Definition 4.6]), then the corresponding matrix sequence  $(A_n)$  has an asymptotic spectrum described by the symbol

$$(6) \quad \omega(x, s) = \frac{k(\phi(x))}{[\phi'(x)]^2} a(s) \quad \text{with} \quad D = (0, 1) \times [-\pi, \pi].$$

For these results, motivated by the use of collocation techniques (see, e.g., [21]) for approximating the solution of one-dimensional and multidimensional boundary value problems, see [35].

In the case of a two-dimensional problem such as (1), the analysis is also quite complete concerning finite difference approximations. For instance, when  $\Omega = (0, 1)^2$  and  $K = I_2$ , using the classical 5 point stencil or the 7 point stencil (in this case there is no difference since  $K_{1,2} = K_{2,1} = 0$ ), we obtain the two-level Toeplitz matrix

$$(7) \quad T_N(b) = T_{n_1}(a) \otimes I_{n_2} + I_{n_1} \otimes T_{n_2}(a),$$

where  $N = (n_1, n_2)$  ( $n_1$  is the number of internal points in the  $x_1$  direction and  $n_2$  is the number of internal points in the  $x_2$  direction),  $n = n_1 n_2$  is the size, and  $b(s_1, s_2) = a(s_1) + a(s_2)$  with  $a(s) = 2 - 2 \cos(s)$ . Also in this case the bivariate stencil represents the nonzero Fourier coefficients of the bivariate generating function  $b$ , and this property remains valid for other stencils. Moreover, according to relation (3), the asymptotic spectrum of  $(T_N(b))_N$  is described by the symbol  $\omega(s_1, s_2) = b(s_1, s_2)$  with  $D = [-\pi, \pi]^2$  (see, e.g., [39]). We observe that the same matrix, with  $n_1 = n_2 = \nu - 1$ , is obtained when employing the  $P_1$  finite element approximation with triangles having the vertices

$$(8) \quad \left( \frac{(j, k)}{\nu}, \frac{(j + \epsilon, k)}{\nu}, \frac{(j, k + \epsilon)}{\nu} \right), \quad \epsilon = \pm 1.$$

More generally, as a consequence of the theory of generalized locally Toeplitz sequences presented in [31, 32], asymptotic spectra can be given for finite difference approximations of (1) for general functions  $K$  and a domain  $\Omega$ , which guarantees the symmetry of the resulting matrix (e.g., a pluri-rectangle that is a connected finite union of rectangles with edges parallel to the main axes; see [36]). For instance, for a 7 point stencil (see the proof of Corollary 1.2(b) below) we know that the resulting matrix sequence has an asymptotic spectrum with symbol

$$(9) \quad \omega(x, s) = \begin{bmatrix} 1 - e^{is_1} \\ 1 - e^{is_2} \end{bmatrix}^* \cdot K(x) \cdot \begin{bmatrix} 1 - e^{is_1} \\ 1 - e^{is_2} \end{bmatrix},$$

with  $D = \Omega \times [-\pi, \pi]^2$ . Notice that if  $\Omega = (0, 1)^2$  and  $K(x) = I_2$ , then the above symbol reduces to that of (7) since

$$\begin{bmatrix} 1 - e^{is_1} \\ 1 - e^{is_2} \end{bmatrix}^* \begin{bmatrix} 1 - e^{is_1} \\ 1 - e^{is_2} \end{bmatrix} = |1 - e^{is_1}|^2 + |1 - e^{is_2}|^2 = a(s_1) + a(s_2) = b(s).$$

Furthermore, for nonequispaced tensor grids obtained as the image under a bijective map  $\phi(x) = (\phi_1(x_1), \phi_2(x_2))^T$  of an equispaced tensor grid, the general structure of the symbol (see [35, 31]) is the natural generalization of (6): denoting by  $\nabla\phi$  the (diagonal) Jacobian of  $\phi(x) = (\phi_1(x_1), \phi_2(x_2))^T$ , we have

$$(10) \quad \omega(x, s) = \begin{bmatrix} 1 - e^{is_1} \\ 1 - e^{is_2} \end{bmatrix}^* \cdot \tilde{K}(x) \cdot \begin{bmatrix} 1 - e^{is_1} \\ 1 - e^{is_2} \end{bmatrix},$$

$$\tilde{K}(x) = \nabla\phi(x)^{-1}K(\phi(x))\nabla\phi(x)^{-T}$$

over  $D = \tilde{\Omega} \times [-\pi, \pi]^2$ ,  $\tilde{\Omega} := \phi^{-1}(\Omega)$ . We notice that (10) is the natural two-dimensional generalization of (6) and that the symbol in (10) reduces to that in (9) if  $\phi_1(x_1) = x_1$  and  $\phi_2(x_2) = x_2$ , i.e., in the case where the grids are uniform.

Finally, recently the above results have been extended to non-Hermitian matrices  $A_n$  occurring, e.g., in the finite difference discretization of PDEs containing lower order difference operators: it has been shown in [16, 18] that, provided that the spectral norm of  $A_n$  is uniformly bounded in  $n$  and that the trace norm of  $S_n = (A_n - A_n^*)/(2i)$ , the skew-Hermitian part of  $A_n$ , grows at most as  $o(n)$ , then the sequence  $(A_n)$  has the same asymptotic spectrum as the sequence  $((A_n + A_n^*)/2)$  obtained from the Hermitian part of  $A_n$ . This result also implies [18, 19] that (9) remains true for more general domains  $\Omega$ , even if one uses different approximation schemes for the boundary conditions.

**1.2. The case of finite element approximations.** Taking into account the results of the previous subsection, the natural question arises of whether similar results on the asymptotic spectrum hold for matrices obtained by applying finite elements to (1). We mentioned already the well-known fact that for the special case  $K = I_2$ ,  $\Omega = (0, 1)^2$  and a uniform triangulation on the square such as (8), the stiffness matrix for  $P_1$  elements is identical to that obtained by finite differences using a 5 point stencil. However, this connection is no longer true in the general case and is not sufficient for us to fully understand the asymptotic properties of stiffness matrices, since for finite elements, for instance, a triangulation does not need to be of tensor form.

Rather than developing a general theory, we will discuss in this paper only the example of an approximation of (1) using  $P_1$  finite elements, together with triangulations  $\mathcal{T}_\nu$  allowing for some a priori mesh refinement. More specifically, in the following we suppose that we have some  $\nu \geq 1$ , some open bounded set  $\tilde{\Omega}$ , and a triangulation  $\mathcal{T}_\nu$  of  $\text{Clos}(\Omega)$  with vertices described by a bijective mapping  $\phi : \text{Clos}(\tilde{\Omega}) \mapsto \text{Clos}(\Omega)$  of the form

$$(11) \quad (j/\nu, k/\nu)^T \in \text{Clos}(\tilde{\Omega}) : \quad P_{j,k} = \phi((j/\nu, k/\nu))$$

and triangles

$$(12) \quad (P_{j,k}, P_{j+\epsilon, k}, P_{j, k+\epsilon}), \quad \epsilon = \pm 1.$$

Such a function  $\phi$  allows us to include also graded triangulations which are suitable if our domain  $\Omega$  has nonconvex vertices (e.g., for L-shaped domains); see Examples 1.3 and 1.4 below. The usual procedure for solving (the variational form of) (1) via  $P_1$  finite elements (see, e.g., [10, 13]) is to consider for  $P_{j,k} \in \Omega$  the hat function  $\psi_{j,k}$  being linear on each of the triangles, taking the value 1 on the vertex  $P_{j,k}$  and 0 on any other vertex (and thus having a support given by the set of the six triangles which

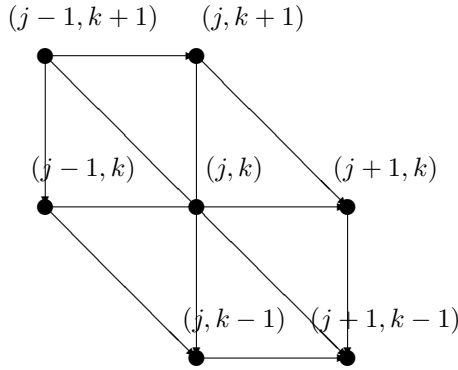


FIG. 1. The vertex  $(j, k)$  and its adjacent vertices for  $P_1$  finite elements.

share the vertex  $P_{j,k}$ ; see Figure 1), and to solve the system of linear equations

$$(13) \quad A_n x_n = b_n, \quad A_n = \left( \int_{\Omega} \nabla \psi_{j,k}(x) K(x) \nabla \psi_{j',k'}(x)^T dx \right)_{P_{j,k}, P_{j',k'} \in \Omega}$$

with a suitable right-hand side  $b_n$  depending on  $f$  and  $g$ . The matrix  $A_n$  is usually referred to as the stiffness matrix. Notice that the same matrix of coefficients but a different right-hand side is obtained if the Dirichlet boundary conditions are partly replaced by Neumann boundary conditions. In what follows, the letter  $n$  will always denote the size of the matrix  $A_n$ , i.e., the number of vertices in  $\Omega$  (which is proportional to  $\nu^2$ ; compare with (18) below).

**THEOREM 1.1.** *Consider the above triangulation  $\mathcal{T}_\nu$  of  $\text{Clos}(\Omega)$  with vertices (11) and triangles (12). We suppose that  $\phi : \text{Clos}(\tilde{\Omega}) \mapsto \text{Clos}(\Omega)$  is bijective,  $m(\tilde{\Omega}) > 0$ , and that there exists an “exceptional” compact set  $\Gamma \subset \text{Clos}(\tilde{\Omega})$  with  $\partial\tilde{\Omega} \subset \Gamma$  and with Lebesgue measure  $m(\Gamma) = 0$  such that  $K \circ \phi$  is continuous in  $\tilde{\Omega} \setminus \Gamma$ , and  $\phi$  is of class  $C^1$  in  $\tilde{\Omega} \setminus \Gamma$ , with nonsingular Jacobian  $\nabla\phi$ . Then an asymptotic spectrum of the stiffness matrices  $A_n$  of (13) for  $\nu \rightarrow \infty$  exists and is given by the formula*

$$\int f d\sigma = \frac{1}{(2\pi)^2} \frac{1}{m(\tilde{\Omega})} \int_{[-\pi, \pi]^2} ds \int_{\tilde{\Omega}} dx f(\omega(x, s)),$$

where

$$\omega(x, s) = \begin{bmatrix} 1 - e^{is_1} \\ 1 - e^{is_2} \end{bmatrix}^* \cdot \tilde{K}(x) \cdot \begin{bmatrix} 1 - e^{is_1} \\ 1 - e^{is_2} \end{bmatrix},$$

$$\tilde{K}(x) = |\det \nabla\phi(x)| \nabla\phi(x)^{-1} K(\phi(x)) \nabla\phi(x)^{-T}.$$

Moreover, this formula for the asymptotic spectrum remains valid if one uses numerical integration for evaluating the entries of  $A_n$ , as long as the quadrature formula has positive weights and integrates constants exactly.

Some consequences of Theorem 1.1 are summarized in the following result.

**COROLLARY 1.2.** *With the notations and assumptions of Theorem 1.1, the following hold:*

- (a) *The sequence of matrices of coefficients  $(A_n)$  has the same asymptotic spectrum as the one obtained by applying  $P_1$  elements on the uniform triangulation (8) to*

the PDE

$$(14) \quad -\nabla(\tilde{K}\nabla^T u) = \tilde{f} \quad \text{on } \tilde{\Omega}, \quad u = \tilde{g} \quad \text{on } \partial\tilde{\Omega}.$$

Moreover, the bilinear form in the weak formulation of problems (1) and (14) are equivalent via variable transformation.

(b) One obtains for  $(A_n)$  the same asymptotic spectrum as that for matrices obtained by applying finite differences based on a 7 point stencil (see Figure 1) to (14). Moreover,  $(A_n)$  is a (reduced) generalized locally Toeplitz sequence in the sense of [32, Definition 3.1], with the symbol  $\omega(x, s)$  of Theorem 1.1.

It is quite instructive to compare the results of Theorem 1.1 and Corollary 1.2 with those of subsection 1.1 for finite difference discretizations. We observe that the symbol in formula (10) and the expression of  $\omega$  in Theorem 1.1 have a similar structure; in particular, we have the same dependency on the domain  $\Omega$  and on the matrix-valued coefficient function  $K$ . Also, the trigonometric polynomials in  $s_1, s_2$  occurring in Theorem 1.1 are the same as those in (10). These polynomials translate the dependency of the asymptotic spectrum on the discretization scheme (5/7 point stencil or  $P_1$  finite elements). The main difference between the two symbols is the dependency on the triangulation described by our function  $\phi$ : in case of finite elements there is an additional factor  $|\det \nabla\phi|$ , leading to a smoother symbol in neighborhoods of points  $x \in \Gamma$  with  $|\det \nabla\phi(x)| = 0$  (corresponding, e.g., to nonconvex edges of  $\Omega$ ; compare with Example 1.3 below), and implying that the finite element matrix of coefficients has fewer eigenvalues of “large” magnitude than the corresponding finite difference matrix of coefficients.

We conclude this section by considering two examples for triangulations  $\mathcal{T}_\nu$  induced by some mapping  $\phi$ .

*Example 1.3.* Suppose that  $\Omega$  is some nonconvex polygon  $\Omega$ , with nonconvex vertices given by  $a_j, j = 1, \dots, p$ , and corresponding inner angles  $\beta_j\pi \in (\pi, 2\pi)$ , and let  $d > 0$  be sufficiently small. Consider the choice  $\Omega = \tilde{\Omega}$  and

$$\phi(x) = \begin{cases} a_j + (x - a_j) \cdot \left(\frac{\|x - a_j\|}{d}\right)^{\beta_j - 1} & \text{for } \|x - a_j\| < d, \\ x & \text{else,} \end{cases}$$

where  $\|\cdot\|$  denotes the Euclidean norm. By construction,  $\phi : \text{Clos}(\Omega) \mapsto \text{Clos}(\Omega)$  is bijective and of class  $\mathcal{C}^1$  in  $\tilde{\Omega} \setminus \Gamma = \{z \in \Omega : \|z - a_j\| \neq d \text{ for } j = 1, 2, \dots, p\}$ . Its Jacobian for  $\|x - a_j\| < d$  is given by

$$\nabla\phi(x) = \frac{\|x - a_j\|^{\beta_j - 1}}{d^{\beta_j - 1}} \left[ I_2 + (\beta_j - 1) \frac{(x - a_j)(x - a_j)^T}{\|x - a_j\|^2} \right],$$

and  $|\det \nabla\phi(x)| = \beta_j (\|x - a_j\|/d)^{2\beta_j - 2}$  tends to 0 for  $x \rightarrow a_j$ . For the inverse of the normalized Jacobian occurring in the symbol of Theorem 1.1 we find

$$\sqrt{|\det(\nabla\phi(x))|} \nabla\phi(x)^{-1} = \sqrt{\beta_j} \left[ I_2 - \left(1 - \frac{1}{\beta_j}\right) \frac{(x - a_j)(x - a_j)^T}{\|x - a_j\|^2} \right].$$

Notice also that  $\|\nabla\phi(x)\|$  is bounded uniformly in  $\tilde{\Omega} \setminus \Gamma$ , implying that the finesse parameter of the triangulation  $\mathcal{T}_\nu$ , i.e., the largest of the diameters of the triangles of this triangulation, is of order  $\mathcal{O}(1/\nu)$ . We finally observe that for triangles where the largest of the distances of the three vertices to  $a_j$  is given by  $t^{\beta_j} \leq d$  have edges with

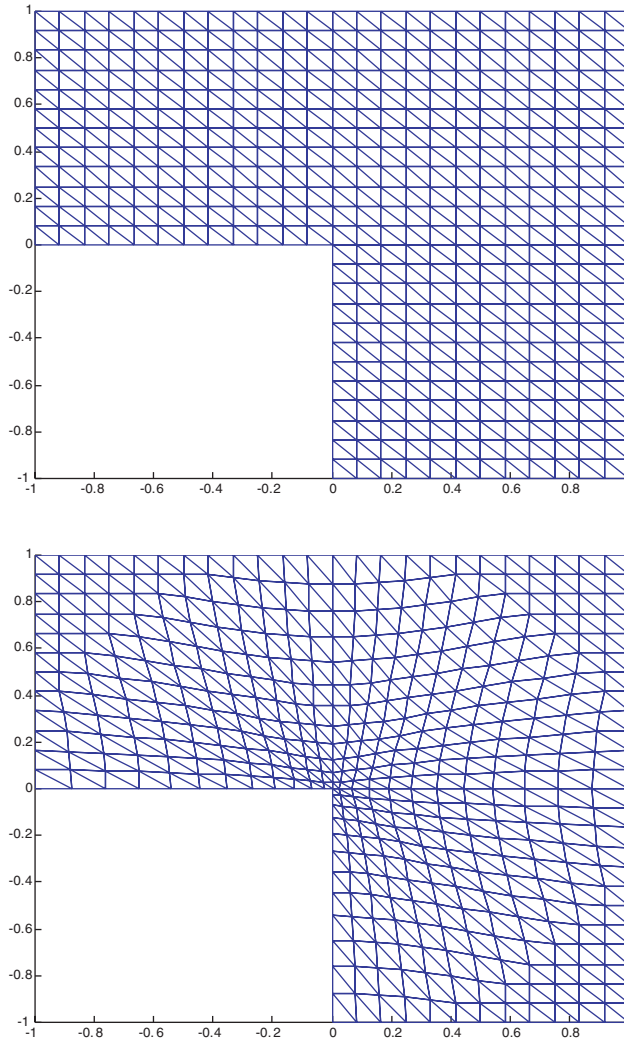


FIG. 2. Triangulation of an L-shape for  $\nu = 12$ . On the top we find the uniform triangulation, and on the bottom its image under the map  $\phi(x) = x \cdot \min\{1, \sqrt{\|x\|}\}$  leading to some grid refinement around the origin.

size of order  $t^{\beta_j-1}/\nu$ : such a mesh refinement based on the grading function  $t \mapsto t^{\beta_j}$  is often used in order to keep the classical finite element error estimate also for singular solutions induced by nonconvex vertices.

*Example 1.4.* A typical example covered by Example 1.3 is a triangulation of an L-shape with vertices  $(0, 0)$ ,  $(-1, 0)$ ,  $(-1, 1)$ ,  $(1, 1)$ ,  $(1, -1)$ ,  $(0, -1)$ , the only nonconvex edge being at the origin  $a_1 = 0$ , with  $\beta_1 := \beta = 3/2$ . Here we can choose  $d = 1$  in Example 1.3, leading to the function  $\phi(x) = x \cdot \min\{1, \|x\|^{\beta-1}\}$ , with the inverse of the normalized Jacobian given by

$$\sqrt{|\det(\nabla\phi(x))|} \nabla\phi(x)^{-1} = \sqrt{\beta} I_2 - \left( \sqrt{\beta} - \frac{1}{\sqrt{\beta}} \right) \frac{xx^T}{\|x\|^2}, \quad \|x\| < 1.$$

In Figure 2 we have drawn both the uniform triangulation and its image under  $\phi$ , leading to some gradation around the origin.



We should notice that in the proof of Theorem 1.1 we do not need any properties of the triangles of  $\mathcal{T}_\nu$  having a nonempty intersection with  $\Gamma$ . Thus Theorem 1.1 remains valid if one uses, for instance, curved elements in order to fit more complicated boundaries.

The remainder of the paper is organized as follows: in section 2 we give the proof of Theorem 1.1 and Corollary 1.2. In section 3 we discuss relations between stiffness matrices for different triangulations, in order to design efficient preconditioning strategies. Finally, in section 4 we draw some conclusions.

**2. Proof.** In what follows we write  $\lambda_1(A_n) \leq \lambda_2(A_n) \leq \dots \leq \lambda_n(A_n)$  for the eigenvalues of some symmetric matrix  $A_n$  of order  $n$ , and  $\mu(A_n) = \frac{1}{n} \sum_{j=1}^n \delta_{\lambda_j(A_n)}$  for the corresponding counting measure. Moreover, we will write  $\mu(A_n) \rightarrow \sigma$  for  $n \rightarrow \infty$  if there is weak-star convergence in the sense of (2), i.e., the matrix sequence  $(A_n)$  has an asymptotic spectrum described by the measure  $\sigma$ .

For proving the above result we make use of the following result on (reduced) generalized locally Toeplitz matrix sequences (see [31, 32]), which we will not cite in its greatest generality: we will focus instead on a subclass of matrix sequences that are (reduced) generalized locally Toeplitz (see [32, Definition 3.1] for the precise definitions in full generality) and also banded and symmetric. Let  $(M_n)$  be a sequence of matrices of size  $n$  and of level  $\gamma \in \mathbb{N}$  defined according to the multi-index rule

$$(15) \quad M_n = (M_{a,a'})_{a,a' \in \nu D \cap \mathbb{Z}^\gamma},$$

$$M_{a,a'} = \frac{1}{(2\pi)^\gamma} \int_{[-\pi, \pi]^\gamma} ds e^{-s^T(a'-a)} \omega\left(\frac{a+a'}{2\nu}, s\right),$$

and corresponding to some open  $D \subset \mathbb{R}^\gamma$ , some integer  $\nu \geq 1$ , and some symbol  $\omega : D \times [-\pi, \pi]^\gamma \rightarrow \mathbb{R}$  with  $\omega(x, s) = \omega(x, -s)$  being a polynomial in  $e^{is}, e^{-is}$  with coefficients continuous in  $x$ . We observe that a matrix  $M_n$  of such a type and level 1 is just an ordinary banded matrix, where succeeding elements on any diagonal vary only slightly (for large  $\nu$  and therefore a fortiori for large  $n$ ) since they are values of some continuous function at arguments differing only by  $1/\nu$  (which tends to zero as  $n = n(\nu)$  tends to infinity). Also, a matrix  $M_n$  of level  $\gamma$  is block banded with blocks being themselves of the same structure as in (15) of level  $\gamma - 1$ . Finally, if the symbol  $\omega(x, s)$  does not depend on  $x$  and  $D = \bigotimes_{j=1}^\gamma (0, \alpha_j)$ , we obtain the classical Toeplitz matrices of level  $\gamma$  and order  $\prod_{j=1}^\gamma [\nu \cdot \alpha_j - 1]$ . A basic result on such symmetric banded (reduced) generalized locally Toeplitz matrix sequences is that they have an asymptotic spectrum given by the following formula [31, 32]:

$$(16) \quad \lim_{n \rightarrow \infty} \mu(M_n) = \sigma, \quad \int f d\sigma = \frac{1}{(2\pi)^\gamma} \frac{1}{m(D)} \int_{[-\pi, \pi]^\gamma} ds \int_D dx f(\omega(x, s)).$$

We will also apply the following statement on the behavior of an asymptotic spectrum under perturbations: the idea relies upon the use of some kind of (matrix) approximation theory for reducing the computation of the symbol of a complicated matrix sequence to the computation of the symbol of simpler matrix sequences (see [29, 31, 32]).

LEMMA 2.1. *Let  $A_n \in \mathbb{C}^{n \times n}$  be symmetric, and suppose that there exist probability measures  $\sigma, \sigma'$  such that, for each  $\epsilon > 0$ , we may write  $A_n = A'_n + A''_n + A'''_n$  with symmetric matrix sequences  $A'_n := A'_n(\epsilon), A''_n := A''_n(\epsilon), A'''_n := A'''_n(\epsilon)$ , where*

$$\limsup_{n \rightarrow \infty} \|A''_n\| \leq \epsilon, \quad \limsup_{n \rightarrow \infty} \frac{\text{rank}(A'''_n)}{n} < \epsilon,$$

and  $(A'_n)_n$  having an asymptotic spectrum  $\mu \leq \epsilon\sigma' + \sigma$ . Then  $(A_n)$  has the asymptotic spectrum  $\sigma$ .

*Proof.* Suppose that the assertion of the lemma is not true. Then by Helley's theorem [25, Theorem 0.1.3] there exists an infinite set of natural numbers  $\mathcal{N}$  such that  $(\mu(A_n))_{n \in \mathcal{N}}$  tends to some probability measure  $\nu$  different from the probability measure  $\sigma$ . By possibly replacing  $A_n$  by  $-A_n$  we may conclude that there exists a  $b \in \mathbb{R}$  with

$$(17) \quad \nu([-\infty, b]) > \sigma([-\infty, b]) = \sigma([-\infty, b]).$$

Write  $r_n = \text{rank}(A''_n)$ . Any  $V \subset \mathbb{C}^n$  can be written as direct sum  $V' \oplus V''$ ,  $V'$  being a subset of the kernel of  $A''_n$ ,  $V''$  being therefore a subset of the image of  $(A''_n)^* = A''_n$ , implying that  $\dim(V') \geq \dim(V) - r_n$ . Consequently, using the Courant min-max principle, we obtain for any  $1 \leq j \leq n - r_n$

$$\begin{aligned} \lambda_j(A'_n) &= \max_{V \subset \mathbb{C}^n, \dim(V)=n+1-j} \min_{y \in V} \frac{y^* A'_n y}{y^* y} \\ &\leq \max_{V \subset \mathbb{C}^n, \dim(V)=n+1-j} \min_{y \in V} \frac{y^* (A'_n + A''_n) y}{y^* y} + \|A''_n\| \\ &\leq \max_{V' \subset \text{Ker}(A''_n), \dim(V') \geq n+1-j-r_n} \min_{y \in V'} \frac{y^* (A'_n + A''_n) y}{y^* y} + \|A''_n\| \\ &\leq \max_{V' \subset \mathbb{C}^n, \dim(V') \geq n+1-j-r_n} \min_{y \in V'} \frac{y^* A_n y}{y^* y} + \|A''_n\| = \lambda_{j+r_n}(A_n) + \|A''_n\|. \end{aligned}$$

Taking into account [25, Theorem 0.1.4], we conclude that

$$\begin{aligned} \nu([-\infty, b]) &\leq \limsup_{n \rightarrow \infty} \mu(A_n)([-\infty, b]) = \limsup_{n \rightarrow \infty} \frac{\#\{j : \lambda_j(A_n) \leq b\}}{n} \\ &\leq \limsup_{n \rightarrow \infty} \frac{r_n + \#\{j > r_n : \lambda_{j-r_n}(A'_n) \leq b + \|A''_n\|\}}{n} \\ &\leq \epsilon + \limsup_{n \rightarrow \infty} \mu(A'_n)([-\infty, b + 2\epsilon]) \leq \epsilon + \sigma([-\infty, b + 2\epsilon]). \end{aligned}$$

For  $\epsilon \rightarrow 0$ , we are left with  $\nu([-\infty, b]) \leq \sigma([-\infty, b])$ , in contradiction with (17). Hence the lemma is shown.  $\square$

The above lemma is essentially contained in original work by Tilli on (one-level) locally Toeplitz sequences [38] and can be considered an evolution of the low-rank, low-norm splittings used by Tyrtyshnikov [39]. A form which is closer to the present approach can be found in [31], where the main role is played by the symbols of the involved matrix sequences. However, in the present version the language and the tools of Lemma 2.1 are a bit different since the results are expressed in terms of measures (recall formulation (2)) rather than symbols (recall formulation (3)).

*Proof of Theorem 1.1.* We start by establishing the formula

$$(18) \quad \lim_{\nu \rightarrow \infty} \frac{n(\nu)}{\nu^2} = m(\tilde{\Omega}), \quad \text{where } n = n(\nu) = \#\left\{ \frac{(j, k)}{\nu} \in \tilde{\Omega} \right\}$$

is the size of the stiffness matrix (13) for the triangulation with parameter  $\nu$ . For  $d > 0$ , denote by  $\Gamma_d := \{y \in \mathbb{R}^2 : \text{dist}(y, \Gamma) \leq d\}$  the closed  $d$ -neighborhood of  $\Gamma$ ,

where we recall that  $\partial\tilde{\Omega} \subset \Gamma$  by assumption on  $\Gamma$ . For any  $\frac{(j,k)}{\nu} \in \tilde{\Omega}$  we find an open square of Lebesgue measure  $1/\nu^2$  being a subset of the  $(2/\nu)$ -neighborhood of  $\tilde{\Omega}$ , any two of such squares having an empty intersection, and thus  $n(\nu)/\nu^2 \leq m(\tilde{\Omega} \cup \Gamma_{2/\nu})$ . On the other hand, the set  $\tilde{\Omega} \setminus \Gamma_{2/\nu}$  is a subset of the union of closed squares of Lebesgue measure  $1/\nu^2$  centered at  $\frac{(j,k)}{\nu} \in \tilde{\Omega}$ , implying that  $n(\nu)/\nu^2 \geq m(\tilde{\Omega} \setminus \Gamma_{2/\nu})$ . Taking into account that  $m(\Gamma_d) \rightarrow m(\Gamma) = 0$  for  $d \rightarrow 0$  by assumption of Theorem 1.1, we arrive at relation (18).

Let  $\epsilon > 0$ . We now choose suitable subsets of  $\tilde{\Omega}$ . Let  $d > 0$  with  $m(\tilde{\Omega} \setminus \Gamma_{3d}) > (1 - \frac{\epsilon}{3}) m(\tilde{\Omega})$ . By compactness of  $\Gamma$ , we may cover  $\Gamma$  with a finite number of open  $\infty$ -neighborhoods  $U_d(x_j) = \{y \in \mathbb{R}^2 : \|y - x_j\|_\infty < d\}$ ,  $j = 1, \dots, p$ , with  $x_j \in \Gamma$ . Defining the pluri-rectangles

$$\tilde{\Omega}' := \tilde{\Omega} \setminus \bigcup_{j=1}^p \text{Clos}(U_{2d}(x_j)), \quad \tilde{\Omega}'' := \tilde{\Omega} \setminus \bigcup_{j=1}^p U_d(x_j),$$

we find that  $\tilde{\Omega} \setminus \Gamma_{3d} \subset \tilde{\Omega}' \subset \tilde{\Omega}'' \subset \tilde{\Omega} \setminus \Gamma$ , with  $\tilde{\Omega}''$  being compact,  $\tilde{\Omega}'$  being open, and

$$(19) \quad \lim_{\nu \rightarrow \infty} \frac{n'(\nu)}{\nu^2} = m(\tilde{\Omega}') \geq \left(1 - \frac{\epsilon}{3}\right) m(\tilde{\Omega}), \quad \text{where } n' = n'(\nu) = \#\left\{\frac{(j,k)}{\nu} \in \tilde{\Omega}'\right\}.$$

Thus, for sufficiently large  $\nu$ ,

$$(20) \quad \frac{n'(\nu)}{n(\nu)} > 1 - \frac{\epsilon}{2}.$$

We are now prepared to introduce a suitable splitting of the stiffness matrix  $A_n$  of (13): we first apply a suitable simultaneous permutation of row and columns such that the first  $n'(\nu)$  rows and columns of  $A_n$  correspond to indices with  $(j,k)/\nu \in \tilde{\Omega}'$ . Then the matrix  $A_n'''$  defined by

$$A_n - A_n''' = \begin{bmatrix} \tilde{A}_n & 0 \\ 0 & 0 \end{bmatrix}, \quad \tilde{A}_n = \left( \int_{\Omega} \nabla \psi_{j,k}(x) K(x) \nabla \psi_{j',k'}(x)^T dx \right)_{(j,k)/\nu, (j',k')/\nu \in \tilde{\Omega}'}$$

is symmetric and has a rank bounded above by twice the difference of the order  $n = n(\nu)$  of  $A_n$  minus the order  $n' = n'(\nu)$  of  $\tilde{A}_n$ . A combination with (20) leads to the relations

$$(21) \quad (A_n''')^* = A_n''', \quad \text{rank}(A_n''') \leq \epsilon n.$$

We want to apply Lemma 2.1 via a splitting  $\tilde{A}_n = \tilde{A}_n' + \tilde{A}_n''$ , and

$$(22) \quad A_n = A_n' + A_n'' + A_n''', \quad A_n' = \begin{bmatrix} \tilde{A}_n' & 0 \\ 0 & 0 \end{bmatrix}, \quad A_n'' = \begin{bmatrix} \tilde{A}_n'' & 0 \\ 0 & 0 \end{bmatrix},$$

where  $\tilde{A}_n''$  will be a symmetric matrix of small norm, and  $\tilde{A}_n'$  symmetric and banded. Moreover,  $(\tilde{A}_n')$  will be (reduced) generalized locally Toeplitz of level 2 in the sense of (15), and thus we know the existence and the explicit form of the asymptotic spectrum of  $(\tilde{A}_n')$  for  $\nu \rightarrow \infty$ .

We make use of the classical assembling procedure of a  $P_1$  finite element matrix  $A_n$ : starting from the zero matrix, the stiffness matrix  $A_n$  is obtained after applying for all triangles  $T$  of the form  $(P_{j,k}, P_{j+\eta,k}, P_{j,k+\eta})$ ,  $\eta = \pm 1$ , the updating formula

(23)

$$A_n \begin{pmatrix} (j, k), (j + \eta, k), (j, k + \eta) \\ (j, k), (j + \eta, k), (j, k + \eta) \end{pmatrix} \leftarrow A_n \begin{pmatrix} (j, k), (j + \eta, k), (j, k + \eta) \\ (j, k), (j + \eta, k), (j, k + \eta) \end{pmatrix} + \frac{1}{2|\det(C^{-1})|} B^T C^{-1} \frac{\int_T K(x) dx}{\int_T dx} C^{-T} B,$$

where the affine mapping  $x \mapsto P_{j,k} + Cx$  brings the points  $(0, 0), (1, 0), (0, 1)$  to  $P_{j,k}, P_{j+\eta,k}, P_{j,k+\eta}$ , respectively, and

$$B = \begin{bmatrix} -1 & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix}.$$

An important observation in our proof is that the updating term in (23) behaves like  $\frac{1}{2}B^T \tilde{K}(\zeta)B$  for some  $\zeta \in \phi^{-1}(T)$  for “most” triangles  $T$ . In order to make this claim more precise in (24) below, we notice that, by construction,  $\tilde{\Omega}''$  is a compact subset of  $\tilde{\Omega} \setminus \Gamma$ , and hence the Jacobian  $\nabla\phi$  of  $\phi$ , its inverse  $\nabla\phi(x)^{-1}$ , and the function  $K \circ \phi$  are uniformly continuous in  $\tilde{\Omega}''$ . Let

$$M := \sup_{x \in \tilde{\Omega}''} \max \left\{ \|\nabla\phi(x)\|, \|\nabla\phi(x)^{-1}\|, \|K(\phi(x))\|, \sqrt{2\epsilon} \right\} \geq 1,$$

and choose  $\nu$  sufficiently large such that a triangle having at least one vertex in  $\tilde{\Omega}'$  is a subset of  $\tilde{\Omega}''$ , and that any of the above functions varies at most by  $\epsilon/(4M^5)$  by choosing two arguments in any triangle that is a subset of  $\tilde{\Omega}''$ . For the matrix  $\tilde{A}_n$  we need to consider only triangles  $T$  having at least one vertex with preimage in  $\tilde{\Omega}'$ . Denoting by  $\tilde{T} \subset \tilde{\Omega}''$  the corresponding triangle with vertices  $\frac{(j,k)}{\nu}, \frac{(j+\eta,k)}{\nu}, \frac{(j,k+\eta)}{\nu}$ , we may conclude with help of the mean value theorem that, for any  $\zeta \in \tilde{T}$ ,

$$\left\| \frac{\int_T K(x) dx}{\int_T dx} - K(\phi(\zeta)) \right\| \leq \frac{\epsilon}{4M^5} \leq M, \quad \left\| \frac{\nu}{\eta} C - \nabla\phi(\zeta) \right\| \leq \frac{\epsilon}{M^5} \leq \frac{1}{2\|\nabla\phi(\zeta)^{-1}\|},$$

and hence

$$\left\| \left(\frac{\nu}{\eta} C\right)^{-1} - \nabla\phi(\zeta)^{-1} \right\| \leq \frac{2\epsilon}{M^3} \leq M, \quad \left\| \det\left(\frac{\nu}{\eta} C\right) - \det(\nabla\phi(\zeta)) \right\| \leq \frac{4\epsilon}{M^4} \leq M.$$

Applying the triangular inequality several times, we obtain after some elementary computations the (quite rough) estimate

$$(24) \quad \max_{\zeta \in \tilde{T}} \left\| \frac{1}{|\det(C^{-1})|} C^{-1} \frac{\int_T K(x) dx}{\int_T dx} C^{-T} - \tilde{K}(\zeta) \right\| \leq 80\epsilon,$$

with  $\tilde{K}$  as in the statement of Theorem 1.1. We remark that the same conclusion holds if instead of exact integration one uses a quadrature formula with positive weights for

TABLE 1

The six adjacent vertices of  $\frac{(j,k)}{\nu} \in \tilde{\Omega}'$  and the corresponding off-diagonal entries of  $\tilde{A}'_n$ : in the first column we find the index  $(j', k')$  of an adjacent vertex, in the second and third column the index of the third vertex of the two triangles giving a nontrivial contribution to the entry in row  $(j, k)$  and column  $(j', k')$  of  $A_n$ , and in the last column the entry of  $\tilde{A}'_n$  at the same position.

$(j', k')$	$(j'', k'')$	$(j''', k''')$	Corresponding entry of $\tilde{A}'_n$
$(j - 1, k)$	$(j, k - 1)$	$(j - 1, k + 1)$	$B_1^T \tilde{K} \left( \frac{(j+j', k+k')}{2\nu} \right) B_2$
$(j, k - 1)$	$(j + 1, k - 1)$	$(j - 1, k)$	$B_1^T \tilde{K} \left( \frac{(j+j', k+k')}{2\nu} \right) B_3$
$(j + 1, k - 1)$	$(j + 1, k)$	$(j, k - 1)$	$B_2^T \tilde{K} \left( \frac{(j+j', k+k')}{2\nu} \right) B_3$
$(j + 1, k)$	$(j, k + 1)$	$(j + 1, k - 1)$	$B_1^T \tilde{K} \left( \frac{(j+j', k+k')}{2\nu} \right) B_2$
$(j, k + 1)$	$(j - 1, k + 1)$	$(j + 1, k)$	$B_1^T \tilde{K} \left( \frac{(j+j', k+k')}{2\nu} \right) B_3$
$(j - 1, k + 1)$	$(j - 1, k)$	$(j, k + 1)$	$B_2^T \tilde{K} \left( \frac{(j+j', k+k')}{2\nu} \right) B_3$

the entries of the stiffness matrix, provided that this quadrature formula integrates constants exactly.

Notice that, in the updating procedure (23), an off-diagonal entry of  $A_n$  is updated twice since a fixed edge of the triangulation is adjacent to two triangles, and a diagonal entry is updated six times since there are six triangles adjacent to a vertex; compare with Figure 1. More precisely, in row labeled  $(j, k)$ , the matrix  $\tilde{A}'_n$  has nonzero off-diagonal entries in columns labeled

$$(j', k') \in \{(j - 1, k + 1), (j, k + 1), (j - 1, k), (j + 1, k), (j, k - 1), (j + 1, k - 1)\},$$

i.e., the indices of adjacent vertices. For instance, for the entry in column  $(j', k') = (j - 1, k)$  we have to consider the two triangles  $T$  with third vertex labeled  $(j'', k'') = (j, k - 1)$ , and  $(j''', k''') = (j - 1, k + 1)$ , respectively, and the corresponding updating quantities can be found at position (1, 2) and (2, 1), respectively, of the symmetric  $3 \times 3$  updating matrix on the right-hand side of (23). Thus, defining the corresponding off-diagonal entry of  $\tilde{A}'_n$  by

$$\begin{aligned} \tilde{A}'_n \left( \begin{matrix} (j', k') \\ (j, k) \end{matrix} \right) &= B_1^T \tilde{K} \left( \frac{1}{2} \left( \frac{(j, k)}{\nu} + \frac{(j', k')}{\nu} \right) \right) B_2 \\ &= (-1, -1) \tilde{K} \left( \frac{(j + j', k + k')}{2\nu} \right) (1, 0)^T, \end{aligned}$$

$B_\ell$  denoting the  $\ell$ th column of  $B$ , we find according to (24) that

$$\left| \tilde{A}'_n \left( \begin{matrix} (j', k') \\ (j, k) \end{matrix} \right) - \tilde{A}_n \left( \begin{matrix} (j', k') \\ (j, k) \end{matrix} \right) \right| \leq 80\epsilon \|B\|^2 = 240\epsilon.$$

The off-diagonal entries of  $\tilde{A}'_n$  for the other five adjacent vertices  $(j', k')$  of  $(j, k)$  are given in Table 1, and each time we obtain the same estimate for the off-diagonal entries of  $\tilde{A}'_n - \tilde{A}_n$ . We define the diagonal entries of  $\tilde{A}'_n$  by

$$\begin{aligned} \tilde{A}'_n \left( \begin{matrix} (j, k) \\ (j, k) \end{matrix} \right) &= \text{trace} \left( B^T \tilde{K} \left( \frac{(j, k)}{\nu} \right) B \right) \\ (25) \quad &= -2 \left( B_1^T \tilde{K} \left( \frac{(j, k)}{\nu} \right) B_2 + B_1^T \tilde{K} \left( \frac{(j, k)}{\nu} \right) B_3 + B_2^T \tilde{K} \left( \frac{(j, k)}{\nu} \right) B_3 \right) \end{aligned}$$

and find according to (24) that

$$\left| \tilde{A}'_n \begin{pmatrix} (j, k) \\ (j, k) \end{pmatrix} - \tilde{A}_n \begin{pmatrix} (j, k) \\ (j, k) \end{pmatrix} \right| \leq 240\epsilon \|B\|^2 = 720\epsilon,$$

and thus, by (22),

$$\|A''_n\| = \|\tilde{A}_n - \tilde{A}'_n\| \leq \sqrt{\|\tilde{A}_n - \tilde{A}'_n\|_1 \|\tilde{A}_n - \tilde{A}'_n\|_\infty} \leq (6 \cdot 240 + 720)\epsilon = 2160\epsilon.$$

It remains to analyze  $\tilde{A}'_n$ . Comparing the definition (15) with the last column of Table 1 and with (25), we see that  $(\tilde{A}'_n)$  is a banded and symmetric generalized locally Toeplitz matrix sequence of level 2 corresponding to the domain  $\tilde{\Omega}'$  and the symbol

$$\begin{aligned} \omega(x, s) &= (2 \cos(s_1) - 2) B_1^T \tilde{K}(x) B_2 + (2 \cos(s_2) - 2) B_1^T \tilde{K}(x) B_3 \\ &\quad + (2 \cos(s_2 - s_1) - 2) B_2^T \tilde{K}(x) B_3 \\ &= 4 \sin^2\left(\frac{s_1}{2}\right) \tilde{K}_{1,1}(x) + 4 \sin^2\left(\frac{s_2}{2}\right) \tilde{K}_{2,2}(x) \\ &\quad + 4 \left[ \sin^2\left(\frac{s_1}{2}\right) + \sin^2\left(\frac{s_2}{2}\right) - \sin^2\left(\frac{s_2 - s_1}{2}\right) \right] \tilde{K}_{1,2}(x), \end{aligned}$$

that is, the same symbol (but a different domain) as in the statement of Theorem 1.1.

Using (16), we may conclude that  $(\mu(\tilde{A}'_n))$  has the limit  $\tilde{\sigma}$ , with

$$\int f d\tilde{\sigma} = \frac{1}{(2\pi)^2} \frac{1}{m(\tilde{\Omega}')} \int_{[-\pi, \pi]^2} ds \int_{\tilde{\Omega}'} dx f(\omega(x, s)).$$

According to (22), for the corresponding counting measures for  $\nu \rightarrow \infty$ , we get using (18), (19),

$$\mu(A'_n) = \frac{n(\nu) - n'(\nu)}{n(\nu)} \cdot \delta_0 + \frac{n'(\nu)}{n(\nu)} \mu(\tilde{A}'_n) \rightarrow \frac{m(\tilde{\Omega}) - m(\tilde{\Omega}')}{m(\tilde{\Omega})} \cdot \delta_0 + \frac{m(\tilde{\Omega}')}{m(\tilde{\Omega})} \tilde{\sigma}$$

and

$$\frac{m(\tilde{\Omega}) - m(\tilde{\Omega}')}{m(\tilde{\Omega})} \cdot \delta_0 + \frac{m(\tilde{\Omega}')}{m(\tilde{\Omega})} \tilde{\sigma} \leq \epsilon \cdot \delta_0 + \frac{m(\tilde{\Omega}')}{m(\tilde{\Omega})} \tilde{\sigma} \leq \epsilon \cdot \delta_0 + \sigma,$$

since  $\tilde{\sigma}$  differs from  $\sigma$  by using a different normalization and a smaller set of integration  $\tilde{\Omega}' \subset \tilde{\Omega}$ . Thus we may apply Lemma 2.1, giving the asymptotic spectrum for  $(A_n)$  as claimed in Theorem 1.1.  $\square$

*Proof of Corollary 1.2.* The first sentence of part (a) follows immediately by applying the formulas of Theorem 1.1 twice. With respect to the second one, consider the variable transformation  $x = \phi(\tilde{x})$  in (1): with  $\tilde{f}(\tilde{x}) = f(\phi(\tilde{x}))$ , we have  $\tilde{\nabla} \tilde{f}(\tilde{x}) = (\nabla f)(\phi(\tilde{x})) \nabla \phi(\tilde{x})$ , and hence

$$\begin{aligned} &\int_{\Omega} (\nabla u)(x) K(x) (\nabla v)(x)^T dx \\ &= \int_{\Omega} (\tilde{\nabla} \tilde{u})(\tilde{x}) \nabla \phi(\tilde{x})^{-1} K(\phi(\tilde{x})) \nabla \phi(\tilde{x})^{-T} (\tilde{\nabla} \tilde{v})(\tilde{x})^T |\det \nabla \phi(\tilde{x})| d\tilde{x} \\ &= \int_{\tilde{\Omega}} (\tilde{\nabla} \tilde{u})(\tilde{x}) \tilde{K}(\tilde{x}) (\tilde{\nabla} \tilde{v})(\tilde{x})^T d\tilde{x}. \end{aligned}$$

For a proof of part (b), we consider

$$y_\nu = (u_{j,k})_{(j,k)/\nu \in \tilde{\Omega}'}, \quad \tilde{u}_{j,k} \approx u \left( \frac{(j,k)}{\nu} \right)$$

and the second order central difference operators using the 7 point stencil of Figure 1,

$$\Delta_1 u_{j,k} = u_{j+1/2,k} - u_{j-1/2,k} \approx \frac{1}{\nu} \frac{\partial}{\partial \tilde{x}_1} u \left( \frac{(j,k)}{\nu} \right),$$

$$\Delta_2 u_{j,k} = u_{j,k+1/2} - u_{j,k-1/2} \approx \frac{1}{\nu} \frac{\partial}{\partial \tilde{x}_2} u \left( \frac{(j,k)}{\nu} \right),$$

$$\Delta_3 u_{j,k} = u_{j+1/2,k-1/2} - u_{j-1/2,k+1/2} \approx \frac{1}{\nu} \left( \frac{\partial}{\partial \tilde{x}_1} - \frac{\partial}{\partial \tilde{x}_2} \right) u \left( \frac{(j,k)}{\nu} \right).$$

Let  $\tilde{\Omega}'$  and  $\tilde{A}'_n$  be as in the preceding proof, and let  $C_n$  be obtained from the matrix  $\tilde{A}'_n$  by replacing the diagonal entries (25) by

$$\begin{aligned} C_n \left( \frac{(j,k)}{(j,k)} \right) &= -B_1^T \left( \tilde{K} \left( \frac{(2j-1, 2k)}{2\nu} \right) + \tilde{K} \left( \frac{(2j+1, 2k)}{2\nu} \right) \right) B_2 \\ &\quad - B_1^T \left( \tilde{K} \left( \frac{(2j, 2k-1)}{2\nu} \right) + \tilde{K} \left( \frac{(2j, 2k+1)}{2\nu} \right) \right) B_3 \\ &\quad - B_2^T \left( \tilde{K} \left( \frac{(2j+1, 2k-1)}{2\nu} \right) + \tilde{K} \left( \frac{(2j-1, 2k+1)}{2\nu} \right) \right) B_3, \end{aligned}$$

and hence  $\|\tilde{A}'_n - C_n\|$  is of order  $\epsilon$ ; compare with (24). For a grid point  $\frac{(j,k)}{\nu} \in \tilde{\Omega}'$  having all its adjacent vertices in  $\tilde{\Omega}'$ , the component of  $C_n y_\nu$  with index  $(j,k)$  can be written as

$$\begin{aligned} &[\tilde{K}_{1,1} + \tilde{K}_{1,2}] \left( \frac{(2j-1, 2k)}{2\nu} \right) (u_{j,k} - u_{j-1,k}) \\ &\quad + [\tilde{K}_{1,1} + \tilde{K}_{1,2}] \left( \frac{(2j+1, 2k)}{2\nu} \right) (u_{j,k} - u_{j+1,k}) \\ &\quad + [\tilde{K}_{2,2} + \tilde{K}_{1,2}] \left( \frac{(2j, 2k-1)}{2\nu} \right) (u_{j,k} - u_{j,k-1}) \\ &\quad + [\tilde{K}_{2,2} + \tilde{K}_{2,1}] \left( \frac{(2j, 2k+1)}{2\nu} \right) (u_{j,k} - u_{j,k+1}) \\ &\quad + \tilde{K}_{1,2} \left( \frac{(2j+1, 2k-1)}{2\nu} \right) (u_{j+1,k-1} - u_{j,k}) \\ &\quad + \tilde{K}_{1,2} \left( \frac{(2j-1, 2k+1)}{2\nu} \right) (u_{j-1,k+1} - u_{j,k}) \\ &= -\Delta_1 [\tilde{K}_{1,1} + \tilde{K}_{1,2}] \Delta_1 u_{j,k} - \Delta_2 [\tilde{K}_{2,2} + \tilde{K}_{1,2}] \Delta_2 u_{j,k} + \Delta_3 \tilde{K}_{1,2} \Delta_3 u_{j,k}. \end{aligned}$$

If some of the vertices  $\frac{(j',k')}{\nu}$  adjacent to  $\frac{(j,k)}{\nu}$  lie outside of  $\tilde{\Omega}'$ , we get a similar expression, where the corresponding values  $u_{j',k'}$  have to be dropped. Therefore the matrix  $C_n$  describes a finite difference discretization in  $\tilde{\Omega}'$  based on the 7 point stencil of Figure 1 for the differential expression

$$\begin{aligned} & -\frac{\partial}{\partial \tilde{x}_1} \left( [\tilde{K}_{1,1} + \tilde{K}_{1,2}] \frac{\partial u}{\partial \tilde{x}_1} \right) - \frac{\partial}{\partial \tilde{x}_2} \left( [\tilde{K}_{2,2} + \tilde{K}_{1,2}] \frac{\partial u}{\partial \tilde{x}_2} \right) \\ & + \left( \frac{\partial}{\partial \tilde{x}_1} - \frac{\partial}{\partial \tilde{x}_2} \right) \left( \tilde{K}_{1,2} \left( \frac{\partial u}{\partial \tilde{x}_1} - \frac{\partial u}{\partial \tilde{x}_2} \right) \right), \end{aligned}$$

coinciding with  $-\nabla(\tilde{K}\nabla u)$ , the differential expression of the PDE of Corollary 1.2(a). Using the same limit considerations as in the proof of Theorem 1.1, the first assertion of Corollary 1.2(b) follows. The second assertion now is a simple consequence of the above relationship between  $A_n$  and the 7 point stencil finite difference matrix and of the fact that every finite difference discretization of second order PDEs leads to (reduced) generalized locally Toeplitz sequences (see [31, 32]).  $\square$

**3. Uniform versus nonuniform triangulations and preconditioning.** Let us briefly recall some classical terminology concerning finite element triangulations. The *finesse parameter* of a triangulation  $\mathcal{T}_\nu$  is the largest among the diameters of the triangles of this triangulation. A family of triangulations  $\mathcal{T}_\nu$  for varying  $\nu$  is called quasi-uniform [2, 20] if the length of the shortest edge in  $\mathcal{T}_\nu$  divided by the finesse parameter of  $\mathcal{T}_\nu$  is bounded below by some constant uniformly in  $\nu$ . The family of triangulations  $\mathcal{T}_\nu$  is called *shape-regular* [9, Definition II.5.1] if the ratio of the diameter divided by the radius of the largest inscribed disk is bounded uniformly for each triangle  $T \in \mathcal{T}_\nu$  and all  $\nu$  (or, equivalently, if all angles are bounded away from zero uniformly in  $\nu$ ).

In the previous sections we have considered a triangulation  $\mathcal{T}_\nu$  of  $\Omega$  being the image under a bijective map  $\phi$  of a uniform triangulation  $\tilde{\mathcal{T}}_\nu$  of  $\tilde{\Omega}$  with stepsize  $1/\nu$ . Denote by  $A_n(K, \mathcal{T}_\nu)$  the corresponding stiffness matrix (13). Since in general the two triangulations  $\mathcal{T}_\nu$  and  $\tilde{\mathcal{T}}_\nu$  lead to stiffness matrices of the same size, we want to discuss in this section in more detail some spectral properties of the matrix  $A_n(I_2, \tilde{\mathcal{T}}_\nu)^{-1} A_n(K, \mathcal{T}_\nu)$  and other related matrices. This analysis is motivated by the task of finding efficient preconditioning strategies for the method of conjugate gradients applied to the stiffness matrix  $A_n(K, \mathcal{T}_\nu)$ . Our uniform triangulation  $(\tilde{\mathcal{T}}_\nu)_\nu$  is trivially both quasi-uniform and shape-regular, while  $(\mathcal{T}_\nu)_\nu$  is not necessarily so. For instance, for the graduated mesh of Example 1.3 we find a finesse parameter  $\geq 1/\nu$ , but the triangle with vertex  $a_j$  has edges of size  $d(1/(d\nu))^{\beta_j}$ , and hence  $(\mathcal{T}_\nu)_\nu$  is not quasi-uniform. In this section we will be particularly interested in the case where  $(\mathcal{T}_\nu)_\nu$  is only shape-regular.

The main results of this section are given in subsection 3.2: in Theorem 3.2 we first relate two stiffness matrices with respect to the partial ordering of Hermitian matrices ( $M_1 \leq M_2$  if  $M_1, M_2$  are Hermitian and  $M_2 - M_1$  is semipositive definite). Subsequently, in Corollary 3.4 we deduce bounds for the smallest and the largest eigenvalue of such preconditioned stiffness matrices, and in Theorem 3.5 we give results on the asymptotic spectrum for such matrices. But first we provide in subsection 3.1 a basic proposition (based on the local analysis of finite element matrices), which is the keystone for proving the results of subsection 3.2.



**3.1. Local domain analysis of the finite element matrices.** In order to better understand the local properties of a stiffness matrix, let us go back to the classical assembling procedure of a  $P_1$  finite element matrix  $A_n$  mentioned already in the proof of Theorem 1.1. Starting from the zero matrix, we have the following updating formulas: any triangle  $T$  of the form  $(P_{j,k}, P_{j+\eta,k}, P_{j,k+\eta})$ ,  $\eta \in \{\pm 1\}$ , gives the contribution

(26)

$$A_n \begin{pmatrix} (j, k), (j + \eta, k), (j, k + \eta) \\ (j, k), (j + \eta, k), (j, k + \eta) \end{pmatrix} \leftarrow A_n \begin{pmatrix} (j, k), (j + \eta, k), (j, k + \eta) \\ (j, k), (j + \eta, k), (j, k + \eta) \end{pmatrix} + \frac{1}{2|\det(C^{-1})|} B^T C^{-1} \frac{\int_T K(x) dx}{\int_T dx} C^{-T} B,$$

where the affine mapping  $x \mapsto P_{j,k} + Cx$  maps the points  $(0, 0), (1, 0), (0, 1)$  to  $P_{j,k}, P_{j+\eta,k}, P_{j,k+\eta}$ , respectively, and

$$B = \begin{bmatrix} -1 & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix}.$$

Suppose that the three points  $(P_{j,k}, P_{j+\eta,k}, P_{j,k+\eta})$  have positive orientation, and define by  $\alpha, \beta, \gamma$ , respectively, the angles of the triangle  $T$  at these vertices. In addition, define  $\Pi$  to be a rotation matrix mapping the half line  $(0, P_{j+\eta,k} - P_{j,k})$  to the half line  $((0, 0), (1, 0))$ ; then

$$\Pi C = \frac{\|P_{j+\eta,k} - P_{j,k}\|}{\sin(\alpha)} \begin{bmatrix} \sin(\gamma) & \sin(\beta) \cos(\alpha) \\ 0 & \sin(\beta) \sin(\alpha) \end{bmatrix},$$

and, in addition,

$$\frac{C^{-1}}{\sqrt{|\det(C^{-1})|}} = \frac{1}{\sqrt{\sin(\alpha) \sin(\beta) \sin(\gamma)}} \begin{bmatrix} \sin(\alpha) \sin(\beta) & -\cos(\alpha) \sin(\beta) \\ 0 & \sin(\gamma) \end{bmatrix} \cdot \Pi.$$

Observe also that  $C^{-1}/\sqrt{|\det(C^{-1})|}$  has the singular values  $\sqrt{\delta_T}$  and  $1/\sqrt{\delta_T}$  and thus a spectral condition number  $\delta_T$ , which can be computed explicitly in terms of the angles of  $T$ :

(27)

$$\delta_T := \text{cond} \left( \frac{C^{-1}}{\sqrt{|\det(C^{-1})|}} \right) = y_T + \sqrt{y_T^2 - 1}, \quad y_T = \frac{\sin^2(\beta) + \sin^2(\gamma)}{2 \sin(\alpha) \sin(\beta) \sin(\gamma)}.$$

Therefore

$$(28) \quad \frac{1}{\delta_T} I_2 \leq \frac{1}{|\det(C^{-1})|} C^{-1} C^{-T} \leq \delta_T I_2.$$

If the three points  $(P_{j,k}, P_{j+\eta,k}, P_{j,k+\eta})$  have negative orientation, then we switch axes; that is, we exchange the role of  $\beta$  and  $\gamma$ , but the conclusions in (27) and (28) are the same. For instance, for a triangle  $T \in \tilde{\mathcal{T}}_\nu$  of a uniform triangulation we get  $\alpha = \pi/2$  and  $\beta = \gamma = \pi/4$ , leading to  $\delta_T = 1$ , but in general  $\delta_T \geq 1$ .

The relation (28) enables us to compare the updating matrices in (26) for different meshes and  $K = I_2$ , and, by a similar argument, for different (pointwise symmetric positive definite) coefficient functions  $K$ .

PROPOSITION 3.1. *With*

$$\kappa_{\min} = \operatorname{ess\,inf}_{x \in T} \lambda_{\min}(K(x)) \geq 0, \quad \kappa_{\max} = \operatorname{ess\,sup}_{x \in T} \lambda_{\max}(K(x)),$$

and  $B, C$  as in (26) we have that

$$\kappa_{\min} \frac{B^T C^{-1} C^{-T} B}{2|\det(C^{-1})|} \leq \frac{1}{2|\det(C^{-1})|} B^T C^{-1} \frac{\int_T K(x) dx}{\int_T dx} C^{-T} B \leq \kappa_{\max} \frac{B^T C^{-1} C^{-T} B}{2|\det(C^{-1})|},$$

and, with  $\delta_T \geq 1$  as in (27),

$$\frac{1}{\delta_T} \frac{B^T B}{2} \leq \frac{B^T C^{-1} C^{-T} B}{2|\det(C^{-1})|} \leq \delta_T \frac{B^T B}{2}.$$

There are many ways of writing the constant  $\delta_T$  of (27). For instance, if  $\beta, \gamma \in (0, \pi/2)$ , we find using the relation  $\alpha + \beta + \gamma = \pi$  that

$$y_T = \frac{\sin^2(\beta) + \sin^2(\gamma)}{\sin^2(\beta) \sin(2\gamma) + \sin^2(\gamma) \sin(2\beta)} \leq \frac{1}{\min(\sin(2\beta), \sin(2\gamma))},$$

which is quite precise if  $\beta$  or  $\gamma$  is small compared to the other two angles. We also have that  $\delta_T$  is uniformly bounded for  $T \in \mathcal{T}_\nu$  for all  $\nu$  if and only if all angles occurring in  $\mathcal{T}_\nu$  are bounded away from zero uniformly in  $\nu$ , i.e.,  $(\mathcal{T}_\nu)_\nu$  is shape-regular. Moreover, there holds

$$\delta_T \leq 2y_T = \frac{b^2 + c^2}{2m(T)} \leq \frac{a + b + c}{2m(T)} \max\{a, b, c\},$$

the expression on the right-hand side being bounded above by the ratio of the diameter of the triangle  $T$  to the radius of the largest disk contained in  $T$ .

For our triangulation  $\mathcal{T}_\nu$  obtained as the image of the uniform triangulation, we also know from the proof of Theorem 1.1 that

$$(29) \quad \frac{C}{\sqrt{|\det(C)|}} \approx \eta \frac{\nabla\phi(\zeta)}{\sqrt{|\det(\nabla\phi(\zeta))|}}, \quad \zeta \in \phi^{-1}(T),$$

and hence

$$\delta := \sup_{\nu} \max_{T \in \mathcal{T}_\nu} \delta_T = \sup_{\nu} \max_{T \in \mathcal{T}_\nu} \operatorname{cond} \left( \frac{C}{\sqrt{|\det(C)|}} \right) \approx \sup_{\zeta \in \bar{\Omega} \setminus \Gamma} \operatorname{cond} \left( \frac{\nabla\phi(\zeta)}{\sqrt{|\det(\nabla\phi(\zeta))|}} \right).$$

This latter quantity turns out to be very simple for the refined triangulations discussed in Examples 1.3 and 1.4, namely  $\delta \approx \beta$ , with  $\beta \in (\pi, 2\pi)$  being the largest inner angle of  $\Omega$ . We should notice that these last arguments are not completely rigorous, since in general relation (29) can be shown to be true only for triangles  $T$  with  $\phi^{-1}(T)$  having a certain distance to  $\Gamma$ . However, there exist similar mesh refinements where the resulting family  $(\mathcal{T}_\nu)_\nu$  is shape-regular and where explicit lower bounds for the angles are known.

### 3.2. Extremal eigenvalues, condition numbers, and preconditioning.

The four statements in this section will have a short proof since they are related to previously known results. For our first statement we have been strongly inspired by similar results for so-called matrix-valued linear and positive operators (LPOs) (see [27, 34]). Here we give a short direct proof.

**THEOREM 3.2.** *Assume that the matrix  $K$  is uniformly elliptic and bounded; i.e., there exist positive constant  $\kappa_{\min}$  and  $\kappa_{\max}$  such that  $\kappa_{\min}I_2 \leq K(x) \leq \kappa_{\max}I_2$  almost everywhere with respect to  $x$  (for instance,  $\kappa_{\min} = \operatorname{ess\,inf}_x \lambda_{\min}(K(x))$ ,  $\kappa_{\max} = \operatorname{ess\,sup}_x \lambda_{\max}(K(x))$ ). Then*

$(A_n(K, \mathcal{T}_\nu))_\nu$  and  $(A_n(I_2, \mathcal{T}_\nu))_\nu$  are uniformly equivalent

$$(30) \quad \text{and more precisely, } \kappa_{\min}A_n(I_2, \mathcal{T}_\nu) \leq A_n(K, \mathcal{T}_\nu) \leq \kappa_{\max}A_n(I_2, \mathcal{T}_\nu),$$

and the same result is true if one replaces  $\mathcal{T}_\nu$  in (30) by  $\tilde{\mathcal{T}}_\nu$ .

Assume that the family of triangulations  $(\mathcal{T}_\nu)_\nu$  is shape-regular, and define

$$\delta := \sup_\nu \max_{T \in \mathcal{T}_\nu} \delta_T < \infty$$

with  $\delta_T$  as in (27). Then

$(A_n(I_2, \mathcal{T}_\nu))_\nu$  and  $(A_n(I_2, \tilde{\mathcal{T}}_\nu))_\nu$  are uniformly equivalent

$$(31) \quad \text{and more precisely } \frac{1}{\delta}A_n(I_2, \tilde{\mathcal{T}}_\nu) \leq A_n(I_2, \mathcal{T}_\nu) \leq \delta A_n(I_2, \tilde{\mathcal{T}}_\nu).$$

*Proof.* The main work for proving statements (30) and (31) has been done already in subsection 3.1: according to (26), the claimed inequalities in (30) are obtained by summing over all triangles  $T \in \mathcal{T}_\nu$  the first inequality of Proposition 3.1. Similarly, relating the triangulations  $\mathcal{T}_\nu$  and  $\tilde{\mathcal{T}}_\nu$  for  $K = I_2$  means that we have to study how the stiffness matrix changes if  $C$  in (26) is replaced by  $I_2$ : the answer is obtained by summing the last inequality of Proposition 3.1 for all triangles (after replacing  $\delta_T$  by  $\delta$ ).  $\square$

The preceding result enables us to give more precise bounds for the smallest and largest eigenvalue of the different stiffness matrices occurring in Theorem 3.2.

**COROLLARY 3.3.** *Assume that the matrix  $K$  is uniformly elliptic and bounded, and that  $(\mathcal{T}_\nu)_\nu$  is shape-regular. Then the largest eigenvalue of  $A_n(K, \mathcal{T}_\nu)$  is uniformly bounded in  $\nu$ , and the smallest behaves like  $1/\nu^2$  for  $\nu \rightarrow \infty$ .*

*In particular, the spectral condition number of  $A_n(K, \mathcal{T}_\nu)$  behaves like  $n$ , the number of vertices of  $\tilde{\mathcal{T}}_\nu$ .*

*Proof.* Since  $\Omega$  is bounded, it is contained in a square with sides of size  $d_{\text{out}}$  and contains a square of size  $d_{\text{in}}$ . Then  $A_n(I_2, \tilde{\mathcal{T}}_\nu)$  contains as submatrix the Toeplitz matrix generated by  $4 - 2\cos(s_1) - 2\cos(s_2)$  of order  $d_{\text{in}}(\nu - 1)^2$ , and, in addition,  $A_n(I_2, \tilde{\mathcal{T}}_\nu)$  is a submatrix of a Toeplitz matrix generated by  $4 - 2\cos(s_1) - 2\cos(s_2)$  of order  $d_{\text{out}}^2\nu^2$  (see [36]). Since the eigenvalues of Toeplitz matrices generated by linear cosine polynomials are explicitly known, it follows that the smallest eigenvalue of  $A_n(I_2, \tilde{\mathcal{T}}_\nu)$  is of order  $1/\nu^2 \sim n^{-1}$ , and its maximal eigenvalue is uniformly bounded by 8, which is also its limit for  $\nu \rightarrow \infty$ . Using, for instance, the well-known representation of extremal eigenvalues of Hermitian matrices in terms of Rayleigh quotients, it follows from Theorem 3.2, by combining (30) and (31), that all three matrices  $A_n(K, \mathcal{T}_\nu)$ ,  $A_n(I_2, \mathcal{T}_\nu)$ , and  $A_n(K, \tilde{\mathcal{T}}_\nu)$  have a smallest eigenvalue of order  $1/\nu^2 \sim n^{-1}$  and a maximal eigenvalue bounded uniformly in  $\nu$ .  $\square$

Corollary 3.3 has been proved in [2, relation (5.102c), p. 235, and pp. 236–238], [9, Lemma V.2.6], [20, p. 61 and Lemma 2.6, p. 233], and [37, Theorem 5.1], under the additional assumption that  $(\mathcal{T}_\nu)_\nu$  is also quasi-uniform. Notice that the proofs given in the above references consist of comparing suitable Sobolev norms, and here the quasi uniformity condition cannot be dropped. The idea contained in subsection 3.1 is to use the updating formulas, i.e., a kind of element-by-element local analysis which is more effective than a global analysis (see, e.g., [1] and the work by Fried [15], where a similar technique has been extensively used).

Let us finally turn to the problem of designing a preconditioner for the CG method applied to the system  $A_n(K, \mathcal{T}_\nu)x_n = b_n$ . We recall that the matrix  $A_n(I_2, \tilde{\mathcal{T}}_\nu)$  corresponding to the uniform triangulation  $\tilde{\mathcal{T}}_\nu$  coincides with that obtained by applying the classical finite difference 5 point stencil to the Poisson problem  $-\Delta u = f$ . Thus solving the system  $A_n(I_2, \tilde{\mathcal{T}}_\nu)y_n = c_n$  can be performed in  $\mathcal{O}(n)$  operations using, e.g., the method of cyclic reductions [11, 12, 14], and thus such a matrix would be a practical preconditioner. Define also the matrix

$$D_n = \text{diag} \left( \left\| \tilde{K} \left( \frac{(j, k)}{\nu} \right) \right\| \right)_{\substack{(j, k) \\ \nu \in \tilde{\Omega}_h}},$$

which again is a practical preconditioner. Then, under the assumptions of Proposition 3.1, the condition number of  $A_n(I_2, \tilde{\mathcal{T}}_\nu)^{-1}A_n(K, \mathcal{T}_\nu)$  and of  $A_n(I_2, \tilde{\mathcal{T}}_\nu)^{-1}D_n^{-1/2} \cdot A_n(K, \mathcal{T}_\nu)D_n^{-1/2}$  can be bounded independently of the stepsize  $1/\nu$  in terms of the smallest angle used in the triangulation of  $\Omega$ , plus possibly the norm and the ellipticity constant of  $K$ . This means that the associated preconditioned CG (PCG) will achieve a fixed precision in  $\mathcal{O}(n)$  operations also in the nonconstant coefficient case with a nonuniform triangulation.

In the following two results we give a complete picture (localization and distribution) of the spectral behavior of preconditioned matrix sequences arising from the use of the above-mentioned preconditioners.

**COROLLARY 3.4.** *Assume that the matrix  $K$  is uniformly elliptic and bounded, i.e., there exist positive constant  $\kappa_{\min}$  and  $\kappa_{\max}$  such that  $\kappa_{\min}I_2 \leq K(x) \leq \kappa_{\max}I_2$  almost everywhere with respect to  $x$  (for instance,  $\kappa_{\min} = \text{essinf}_x \lambda_{\min}(K(x))$ ,  $\kappa_{\max} = \text{esssup}_x \lambda_{\max}(K(x))$ ). Then*

$$(32) \quad \text{the eigenvalues of } A_n(I_2, \mathcal{T}_\nu)^{-1}A_n(K, \mathcal{T}_\nu) \text{ belong to } [\kappa_{\min}, \kappa_{\max}],$$

and the same result is true if one replaces  $\mathcal{T}_\nu$  in (32) by  $\tilde{\mathcal{T}}_\nu$ .

Assume also that the family of triangulations  $(\mathcal{T}_\nu)_\nu$  is shape-regular such that  $\delta := \sup_\nu \max_{T \in \mathcal{T}_\nu} \delta_T < \infty$  with  $\delta_T$  as in (27). Then

$$(33) \quad \text{the eigenvalues of } A_n(I_2, \tilde{\mathcal{T}}_\nu)^{-1}A_n(I_2, \mathcal{T}_\nu) \text{ belong to } [1/\delta, \delta];$$

$$(34) \quad \text{the eigenvalues of } A_n(I_2, \tilde{\mathcal{T}}_\nu)^{-1}A_n(K, \mathcal{T}_\nu) \text{ belong to } [\kappa_{\min}/\delta, \kappa_{\max}\delta].$$

*Proof.* Statements (32) and (33) follow from the corresponding statements (30) and (31) in Theorem 3.2 and the fact that, for Hermitian positive definite  $X, Y$ , we have for the spectrum  $\Lambda(Y^{-1}X)$  the localization

$$\Lambda(Y^{-1}X) \subset \left\{ \frac{u^*Xu}{u^*Yu} : u \neq 0 \right\}.$$

The claim (34) follows from (30) and (31) by rewriting the Rayleigh quotient as

$$\frac{u^* A_n(K, \mathcal{T}_\nu) u}{u^* A_n(I_2, \tilde{\mathcal{T}}_\nu) u} = \frac{u^* A_n(K, \mathcal{T}_\nu) u}{u^* A_n(I_2, \mathcal{T}_\nu) u} \frac{u^* A_n(I_2, \mathcal{T}_\nu) u}{u^* A_n(I_2, \tilde{\mathcal{T}}_\nu) u}. \quad \square$$

**THEOREM 3.5.** *Assume that the matrix  $K$  is uniformly elliptic in the sense of Corollary 3.4. Consider the preconditioned sequences  $(Y_n^{-1} X_n)$  with*

$$[Y_n, X_n] \in \{[A_n(I_2, \tilde{\mathcal{T}}_\nu), A_n(K, \tilde{\mathcal{T}}_\nu)], [A_n(I_2, \mathcal{T}_\nu), A_n(K, \mathcal{T}_\nu)], [A_n(I_2, \tilde{\mathcal{T}}_\nu), A_n(I_2, \mathcal{T}_\nu)], [A_n(I_2, \tilde{\mathcal{T}}_\nu), A_n(K, \mathcal{T}_\nu)]\}.$$

Then, calling  $\omega_X$  the symbol of  $(X_n)$  and calling  $\omega_Y$  the symbol of  $(Y_n)$ , we have that the asymptotic spectrum of  $(Y_n^{-1} X_n)$  is given by  $\omega_X / \omega_Y$ .

*Proof.* It is enough to observe that all the involved matrix sequences are such that both  $X_n$  and  $Y_n$  come from the same matrix-valued LPO for which the distribution is known (see Theorem 1.1) and is sparsely vanishing (i.e., the symbol vanishes in a set of zero Lebesgue measure). The conclusion follows from the general theory of LPOs as in Theorem 2.9 of [28] (compare also Theorem 4.6 in [33] and Theorem 3.7 in [26]).  $\square$

With the notation of the above theorem, we remark that the same result could be proved for the matrices  $[Y_n, X_n] = [D_n^{1/2} A_n(I_2, \tilde{\mathcal{T}}_\nu) D_n^{1/2}, A_n(K, \mathcal{T}_\nu)]$ . Indeed  $D_n^{1/2}$ ,  $A_n(I_2, \tilde{\mathcal{T}}_\nu)$ , and  $A_n(K, \mathcal{T}_\nu)$  are all (reduced) generalized locally Toeplitz sequences with sparsely vanishing symbols (i.e., zero on at most a set of zero Lebesgue measure): for  $D_n$  the statement is trivial since the matrix is diagonal, while for the remaining two matrix sequences this has been proved in Corollary 1.2. Then our claims follow from the fact that, if the symbols are all sparsely vanishing and sparsely unbounded (the inverse of a sparsely vanishing), then the operation  $X_n \odot Y_n$  also gives a sequence in the generalized locally Toeplitz class, with asymptotic spectrum described by the symbol  $\omega_X \odot \omega_Y$ ; this has been shown in [31, Theorem 5.8] for  $\odot$  being multiplication, in the same paper for  $\odot$  being addition or subtraction, and is known to be true also for inversion, that is, for the sequence  $(Y_n^{-1} X_n)$  (see [32, Theorems 2.2 and 3.2]).

In order to illustrate Theorem 3.5 and its link with Theorem 1.1, we mention more explicitly the example that the sequence of matrices  $(A_n(I_2, \tilde{\mathcal{T}}_\nu)^{-1} A_n(K, \mathcal{T}_\nu))$  for  $\nu \rightarrow \infty$  has an asymptotic spectrum described by the measure  $\sigma$ , with

$$\int f d\sigma = \frac{1}{(2\pi)^2} \frac{1}{m(\tilde{\Omega})} \int_{[-\pi, \pi]^2} ds \int_{\tilde{\Omega}} dx f(\omega(x, s)),$$

$\tilde{K}(x) = |\det \nabla \phi(x)| \nabla \phi(x)^{-1} K(\phi(x)) \nabla \phi(x)^{-T}$  as before,

$$\omega(x, s) = \frac{\omega_X(x, s)}{\omega_Y(x, s)} = \frac{\begin{bmatrix} 1 - e^{is_1} \\ 1 - e^{is_2} \end{bmatrix}^* \cdot \tilde{K}(x) \cdot \begin{bmatrix} 1 - e^{is_1} \\ 1 - e^{is_2} \end{bmatrix}}{\begin{bmatrix} 1 - e^{is_1} \\ 1 - e^{is_2} \end{bmatrix}^* \cdot \begin{bmatrix} 1 - e^{is_1} \\ 1 - e^{is_2} \end{bmatrix}},$$

and with  $\omega_X(x, s)$ ,  $\omega_Y(x, s)$  according to the notation of Theorem 3.5.

In particular (compare with (34)), the most important part of its eigenvalues lies in the interval

$$[\kappa_{\min}, \kappa_{\max}] = \left[ \operatorname{ess\,inf}_{x \in \tilde{\Omega}} \lambda_{\min}(\tilde{K}(x)), \operatorname{ess\,sup}_{x \in \tilde{\Omega}} \lambda_{\max}(\tilde{K}(x)) \right].$$

**4. Concluding remarks.** We have shown the existence of an asymptotic spectrum for the sequence of stiffness matrices, which occur in the  $P_1$  finite element approximation of the two-dimensional model problem (1) with an a priori mesh refinement and varying stepsizes. The underlying symbol  $\omega$  of this asymptotic spectrum, given in Theorem 1.1, depends not only on the domain and the coefficient functions of the PDE, but also on the particular  $P_1$  approximation scheme (via the dependency on the Fourier variable  $s$ ) and the map  $\phi$  which describes our mesh refinement. We expect, by analogy with the finite difference case (see [31]), that Theorem 1.1 holds also for other finite elements if one adapts the choice of the trigonometric polynomials in  $s$ . It is probably also possible to extend our results to higher dimensions and to other elliptic PDEs, and probably we need only quite weak regularity assumptions on the involved domain and the involved coefficient functions, as in the finite difference case (see [38, 31, 32]). On the other hand, the graded meshes used in modern solvers (especially those generated by a posteriori mesh refinements) in general are not topologically equivalent to the meshes considered in this paper. Notice that, for proving asymptotic spectral results of global type, it is sufficient that the graded meshes are equivalent to an approximation of our meshes (see [35]). These issues should be investigated in more detail in future works, in order to widen the practical impact of our findings.

In the second part of the paper we have analyzed the spectral behavior of some preconditioned finite element matrix sequences in terms of localization, extremal, and, especially, distributional spectral results. The analysis could be used for deducing more precise bounds on the (P)CG convergence, in view of the results in [4, 5, 6]: the related specific study and the related numerical experiments will be part of a subsequent work.

Beside the locally Toeplitz idea, we have used in section 3.1 another purely linear algebra tool, namely the local domain analysis: it consists of decomposing complicated matrix structures in linear combinations of nonnegative definite dyads or low-rank matrices for which the (spectral) analysis is very simple, and then combining these results to deduce properties of the original matrix. (For finite elements see [1] and the beautiful and rich paper by Fried [15, e.g., (47)]; for finite differences compare with [33, section 3.5], [7, Theorem 3.7]; and for general matrices see [26].) We mention that this simple tool is especially useful for preconditioning analysis and for the analysis of extremal eigenvalues asymptotics. As a byproduct we have deduced in Corollary 3.4 that the finite element matrix sequence with uniform triangulation and the nonuniform one (not necessarily verifying the quasi uniformity) are spectrally equivalent. Thus a simpler (projected) two-level Toeplitz structure associated with the uniform triangulation can be employed as preconditioner requiring a constant number of iterations independently of the size of the problem.

## REFERENCES

- [1] T. APEL, *Anisotropic Finite Elements: Local Estimates and Applications*, Adv. Numer. Math., Teubner, Stuttgart, Germany, 1999.
- [2] O. AXELSSON AND V. BARKER, *Finite Element Solution of Boundary Value Problems, Theory and Computation*, Academic Press, New York, 1984.
- [3] O. AXELSSON AND G. LINDSKOG, *On the rate of convergence of the preconditioned conjugate gradient method*, Numer. Math., 48 (1986), pp. 499–523.
- [4] B. BECKERMANN AND A. B. J. KUIJLAARS, *Superlinear convergence of conjugate gradients*, SIAM J. Numer. Anal., 39 (2001), pp. 300–329.
- [5] B. BECKERMANN AND A. B. J. KUIJLAARS, *On the sharpness of an asymptotic error estimate for conjugate gradients*, BIT, 41 (2001), pp. 856–867.

- [6] B. BECKERMANN AND A. B. J. KUIJLAARS, *Superlinear CG convergence for special right-hand sides*, Electron. Trans. Numer. Anal., 14 (2002), pp. 1–19.
- [7] D. BERTACCINI, G. GOLUB, S. SERRA CAPIZZANO, AND C. TABLINO POSSIO, *Preconditioned HSS method for the solution of non-Hermitian positive definite linear systems and applications to the discrete convection-diffusion equation*, Numer. Math., 99 (2005), pp. 441–484.
- [8] A. BÖTTCHER AND B. SILBERMANN, *Introduction to Large Truncated Toeplitz Matrices*, Springer, New York, 1999.
- [9] D. BRAESS, *Finite Elements: Theory, Fast Solvers, and Applications in Solid Mechanics*, Cambridge University Press, Cambridge, UK, 2001.
- [10] F. BREZZI AND M. FORTIN, *Mixed and Hybrid Finite Element Methods*, Springer, New York, 1991.
- [11] B. L. BUZBEE, G. H. GOLUB, AND C. W. NIELSON, *On direct methods for solving Poisson's equations*, SIAM J. Numer. Anal., 7 (1970), pp. 627–656.
- [12] B. L. BUZBEE, F. W. DORR, J. A. GEORGE, AND G. H. GOLUB, *The direct solutions of the discrete Poisson equation on irregular regions*, SIAM J. Numer. Anal., 8 (1971), pp. 722–736.
- [13] P. CIARLET, *The Finite Element Method for Elliptic Problems*, North-Holland, Amsterdam, 1978.
- [14] F. W. DORR, *The direct solution of the discrete Poisson equation on a rectangle*, SIAM Rev., 12 (1970), pp. 248–263.
- [15] I. FRIED, *Bounds on the spectral and maximum norms of the finite element stiffness, flexibility and mass matrices*, Internat. J. Solids Structures, 9 (1973), pp. 1013–1034.
- [16] L. GOLINSKII AND S. SERRA CAPIZZANO, *The asymptotic properties of the spectrum of non symmetrically perturbed Jacobi matrix sequences*, J. Approx. Theory, 144 (2007), pp. 84–102.
- [17] U. GREXANDER AND G. SZEGÖ, *Toeplitz Forms and Their Applications*, 2nd ed., Chelsea, New York, 1984.
- [18] S. HOLMGREN, S. SERRA CAPIZZANO, AND P. SUNDQVIST, *Can one hear the composition of a drum?* Mediterranean J. Math., 3 (2006), pp. 227–249.
- [19] S. HOLMGREN, S. SERRA CAPIZZANO, AND P. SUNDQVIST, *On the asymptotic spectrum of (non symmetric) finite difference matrix sequences*, 2006, in preparation.
- [20] C. JOHNSON, *Numerical Solutions of Partial Differential Equations by the Finite Elements Methods*, Cambridge University Press, Cambridge, UK, 1988.
- [21] S. D. KIM AND S. V. PARTER, *Preconditioning Chebyshev spectral collocation by finite-differences operators*, SIAM J. Numer. Anal., 34 (1997), pp. 939–958.
- [22] A. B. J. KUIJLAARS AND S. SERRA CAPIZZANO, *Asymptotic zero distribution of orthogonal polynomials with discontinuously varying recurrence coefficients*, J. Approx. Theory, 113 (2001), pp. 142–155.
- [23] S. PARTER, *On the eigenvalues of certain generalizations of Toeplitz matrices*, Arch. Ration. Math. Mech., 3 (1962), pp. 244–257.
- [24] Y. SAAD, *Iterative Methods for Sparse Linear Systems*, PWS Publishing, Boston, 1996.
- [25] E. B. SAFF AND V. TOTIK, *Logarithmic Potentials with External Fields*, Springer, Berlin, 1997.
- [26] S. SERRA CAPIZZANO, *Locally X matrices, spectral distributions, preconditioning, and applications*, SIAM J. Matrix Anal. Appl., 21 (2000), pp. 1354–1388.
- [27] S. SERRA CAPIZZANO, *Some theorems on linear positive operators and functionals and their applications*, Comput. Math. Appl., 39 (2000), pp. 139–167.
- [28] S. SERRA CAPIZZANO, *A note on the asymptotic spectra of finite difference discretizations of second order elliptic partial differential equations*, Asian J. Math., 4 (2000), pp. 499–514.
- [29] S. SERRA CAPIZZANO, *Distribution results on the algebra generated by Toeplitz sequences: A finite dimensional approach*, Linear Algebra Appl., 328 (2001), pp. 121–130.
- [30] S. SERRA CAPIZZANO, *Spectral behaviour of matrix sequences and discretized boundary value problems*, Linear Algebra Appl., 337 (2001), pp. 37–78.
- [31] S. SERRA CAPIZZANO, *Generalized locally Toeplitz sequences: Spectral analysis and applications to discretized partial differential equations*, Linear Algebra Appl., 366 (2003), pp. 371–402.
- [32] S. SERRA CAPIZZANO, *GLT sequences as a generalized Fourier analysis and applications*, Linear Algebra Appl., 419 (2006), pp. 180–233.
- [33] S. SERRA CAPIZZANO AND C. TABLINO POSSIO, *Spectral and structural analysis of high precision finite difference matrices for elliptic operators*, Linear Algebra Appl., 293 (1999), pp. 85–131.
- [34] S. SERRA CAPIZZANO AND C. TABLINO POSSIO, *Finite element matrix-sequences: The case of rectangular domains*, Numer. Algorithms, 28 (2001), pp. 309–327.
- [35] S. SERRA CAPIZZANO AND C. TABLINO POSSIO, *Analysis of preconditioning strategies for collocation linear systems*, Linear Algebra Appl., 369 (2003), pp. 41–75.
- [36] S. SERRA CAPIZZANO AND C. TABLINO POSSIO, *Superlinear preconditioners for finite differences linear systems*, SIAM J. Matrix Anal. Appl., 25 (2003), pp. 152–164.

- [37] G. STRANG AND G. FIX, *An Analysis of the Finite Element Method*, Prentice-Hall, Englewood Cliffs, NJ, 1973.
- [38] P. TILLI, *Locally Toeplitz sequences: Spectral properties and applications*, *Linear Algebra Appl.*, 278 (1998), pp. 91–120.
- [39] E. E. TYRTYSHNIKOV, *A uniform approach to some old and new theorems on distribution and clustering*, *Linear Algebra Appl.*, 232 (1996), pp. 1–43.
- [40] E. E. TYRTYSHNIKOV, *A matrix view on the root distribution for orthogonal polynomials*, in *Structured Matrices: Recent Developments in Theory and Computation*, D. Bini, E. Tyrtyshnikov, and P. Yalamov, eds., Nova Science, New York, 2001, pp. 149–156.
- [41] A. VAN DER SLUIS AND H. A. VAN DER VORST, *The rate of convergence of conjugate gradients*, *Numer. Math.*, 48 (1986), pp. 543–560.