

# STABILITY ANALYSIS OF A SIMPLE DISCRETIZATION METHOD FOR A CLASS OF STRONGLY SINGULAR INTEGRAL EQUATIONS

MARTIN COSTABEL, MONIQUE DAUGE, AND KHADIJEH NEDAIASL

ABSTRACT. Motivated by the discrete dipole approximation (DDA) for the scattering of electromagnetic waves by a dielectric obstacle that can be considered as a simple discretization of a Lippmann-Schwinger style volume integral equation for time-harmonic Maxwell equations, we analyze an analogous discretization of convolution operators with strongly singular kernels.

For a class of kernel functions that includes the finite Hilbert transformation in 1D and the principal part of the Maxwell volume integral operator used for DDA in dimensions 2 and 3, we show that the method, which does not fit into known frameworks of projection methods, can nevertheless be considered as a finite section method for an infinite block Toeplitz matrix. The symbol of this matrix is given by a Fourier series that does not converge absolutely. We use Ewald's method to obtain an exponentially fast convergent series representation of this symbol and show that it is a bounded function, thereby allowing to describe the spectrum and the numerical range of the matrix.

It turns out that this numerical range includes the numerical range of the integral operator, but that it is in some cases strictly larger. In these cases the discretization method does not provide a spectrally correct approximation, and while it is stable for a large range of the spectral parameter  $\lambda$ , there are values of  $\lambda$  for which the singular integral equation is well posed, but the discretization method is unstable.

## CONTENTS

1. Introduction	2
1.1. Motivation	2
1.2. The Discrete Dipole Approximation	4
1.3. Outline of the paper	5
1.4. Notation for Fourier transforms and Fourier series	6
1.5. Kernels and their symbols	7
1.6. Delta-delta discretization	11
2. The discrete system	11
2.1. Toeplitz structure	12

---

1991 *Mathematics Subject Classification.* 65R20, 45E10, 47A12, 65B10, 78M99.

*Key words and phrases.* volume integral equation, strongly singular kernel, delta-delta discretization, discrete dipole approximation, numerical stability.

2.2. Ewald method	12
2.3. An integral representation	15
2.4. Matrix-valued kernels	17
3. Examples	18
3.1. Example 1. Dimension $d = 1$ . Finite Hilbert transformation	18
3.2. Example 2. Dimension $d = 2$ , kernel $x_1 x_2  x ^{-4}$	20
3.3. Example 3. Dimension $d = 2$ , kernel $(x_1^2 - x_2^2)  x ^{-4}$	23
3.4. Example 4. Dimension $d = 2$ , kernel $(x_1 + i x_2)^2  x ^{-4}$	27
3.5. Example 5. Dimension $d \geq 2$ . Volume Integral Equation for the Quasi-static Maxwell system	28
Acknowledgment	33
References	33

## 1. INTRODUCTION

### 1.1. Motivation.

Introduced almost 50 years ago by Purcell and Pennypacker [13], the Discrete Dipole Approximation (DDA) is a classical numerical method in computational electromagnetics that is the subject of a vast and still rapidly growing literature (see the surveys [18, 3]), but is virtually unknown in the mathematical community. It can be considered as a numerical approximation scheme for a strongly singular volume integral equation that, however, is too simple to fit into any known framework for standard approximation schemes for such equations (Galerkin, collocation or Nyström methods etc). In particular, to the authors' knowledge, there does not exist any error estimate or convergence proof for this method.

In the paper [19], estimates for a consistency error are derived, and it is observed that to complete the convergence analysis, a uniform estimate for the inverse of the matrix of the linear system (stability estimate) would be needed. In the present paper, we prove first results on the way to such stability estimates for the DDA and related numerical schemes. The class of singular integral equations considered here includes the quasi-static case (i.e. zero frequency limit) of the Maxwell volume integral equation that describes the scattering of electromagnetic waves by a penetrable dielectric body in the case of constant electric permittivity. Further stability results for the non-zero frequency case will be the subject of a forthcoming paper.

Because of the simplicity of the class of operators considered here (convolution operators with kernels positively homogeneous of degree  $-d$  on a bounded domain  $\Omega \subset \mathbb{R}^d$ ), we are able to obtain rather sharp results on the region of stability, by estimating the numerical range of the discretized operator in comparison with the numerical range of the integral operator. It turns

out that for some operators, including the quasi-static Maxwell case in dimension  $d \geq 2$ , the stability region is smaller than what one would naïvely expect. This corresponds to the fact that the eigenvalues of the system matrix, as the mesh-width of the discretization tends to zero, accumulate on a set that is strictly larger than the convex hull of the essential spectrum of the integral operator.

In the paper [14], motivated by the convergence analysis of iterative solutions of the resulting large linear systems, the essential spectrum of the Maxwell volume integral operator was studied for the case of scattering by a dielectric ball in  $\mathbb{R}^3$ . This is a subset of the segment in the complex plane that corresponds to the essential numerical range of the integral operator. It is now known (see [4, 5]) that the same form of the essential spectrum is valid for more general bounded Lipschitz domains. In [14], results of some numerical experiments are then shown that seem to indicate that the eigenvalues of the system matrices accumulate either at isolated points, corresponding to eigenvalues of the integral operator and hence to eigenvalues or resonances of the scattering problem, or at the points of the segment that is spanned by the essential spectrum of the integral operator. Looking closer at Figures 4.2–4.4 of [14], one can detect an “overshoot”, namely that the observed segment of accumulation points is actually larger than the span of the essential spectrum.

In the paper [20], there is a discussion of the spectrum of the system matrices of the DDA scheme for the quasi-static Maxwell equations, motivated by the numerical modeling of the scattering of light by dust particles whose size is small with respect to the wavelength of the light (“Rayleigh particles”). Based on extensive experience with numerical computations using the DDA code ADDA, the authors are convinced that the DDA provides a faithful approximation of the solution of the volume integral equation in the sense that, among other things, the spectral measure of the DDA system matrices converges to the spectral measure of the volume integral operator. They provide plots of the spectral density of these matrices, including a zoom on a neighborhood of the lower end of the spectrum, see graph (a) in [20, FIG. 8]. There one can clearly see that there is an overshoot, namely a part of the spectrum below zero, and that its negative minimum does not disappear as the number of dipoles grows, but rather seems to converge to some number around  $-0.09$ . In a subsequent paper [15], the authors detect this “spill-out” of the spectrum of the DDA system matrices and relate it to an explosion of the needed iterations in an iterative solution method that they observed for large refractive indices. They study the behavior of this overshoot for anisotropic meshes, where it becomes larger, and for some recently introduced improvements of the DDA, where it seems to disappear.

For the quasi-static Maxwell case we prove below (see Proposition 3.15 and (3.51)) that for the classical DDA on a cubic grid such an overshoot indeed exists and that it amounts to an almost 20% increase of the length of the segment spanned by the essential spectrum.

This somewhat unexpected result implies that the simple discretization scheme of the DDA does not provide a spectrally correct approximation of the strongly singular volume integral operator. The additional observation, supported by numerical experience, that this concerns only a small neighborhood of the essential spectrum or perhaps even only of the endpoints of this spectrum, whereas discrete eigenvalues and large parts of the spectral density nevertheless are correctly approximated, still awaits a precise description and proof.

It also implies that the DDA scheme is actually unstable in high-contrast situations, namely if the relative permittivity is very small (smaller than  $\sim 0.093$ ) or very large (larger than  $\sim 11.8$ ). We prove this here for the zero-frequency limit, but expect that it is also true for non-zero frequencies.

## 1.2. The Discrete Dipole Approximation.

As its name indicates, the DDA (sometimes called Coupled Dipole Approximation) can be considered as an approximation of a dielectric continuum described by Maxwell's equations by a different physical system consisting of a finite number of dipoles that are characterized by their polarizability, interacting via electromagnetic fields.

The same mathematical system can be obtained by a procedure more amenable to arguments of numerical analysis, namely by transforming the Maxwell equations for the original dielectric continuum into an equivalent Lippmann-Schwinger style volume integral equation and then discretizing this integral equation by a simple delta-delta approximation on a regular grid  $\{x_n \mid n \in \mathbb{Z}^d\}$  of meshwidth  $h > 0$ .

Thus a linear integral equation on a bounded domain  $\Omega \subset \mathbb{R}^d$

$$(1.1) \quad \lambda u(x) - \int_{\Omega} K(x, y)u(y)dy = f(x) \quad (x \in \Omega)$$

will be approximated by the finite dimensional linear system

$$(1.2) \quad \lambda u_m - \sum_{x_n \in \Omega, n \neq m} h^d K(x_m, x_n)u_n = f(x_m) \quad (x_m \in \Omega).$$

We omit the diagonal term  $m = n$ , because we shall have to do with singular kernels. Apart from this, (1.2) looks like a Galerkin method with Dirac deltas as trial and test functions.

Let us briefly describe the construction of the volume integral equation. A more detailed derivation can be found in [10] and, with special emphasis on the two-dimensional situation, in [5]. We write the time harmonic Maxwell equations with normalized frequency  $\kappa \in \mathbb{C}$  as a second order system for the electric field  $u$ .

$$(1.3) \quad \text{curl curl } u - \kappa^2 \epsilon u = i\kappa J.$$

Here it is assumed that the magnetic permeability is constant (normalized to 1) in the whole space. If one further assumes that the permittivity  $\epsilon$  is equal to 1 outside of a bounded domain and the source current  $J$  has compact support, one can write this as a perturbation of the free-space situation

$$(1.4) \quad \text{curl curl } u - \kappa^2 u = -\kappa^2(1 - \epsilon)u + i\kappa J.$$

Here the right hand side has compact support, and therefore convolution with the outgoing fundamental solution  $g_{\kappa}$  of the Helmholtz equation and application of the operator  $\nabla \text{div} + \kappa^2$  leads to the volume integral equation in distributional form

$$(1.5) \quad u = -(\nabla \text{div} + \kappa^2)g_{\kappa} \star (1 - \epsilon)u + u^{\text{inc}}.$$

Here the incoming field  $u^{\text{inc}}$  combines the field generated by the current density with possible sourceless full space solutions of Maxwell's equations (plane waves etc.)

Equation (1.5) can be considered in any dimension  $d \geq 2$ , but only  $d = 2$  and  $d = 3$  are relevant for electrodynamics. The equation can be written in the form of a second kind strongly singular integral equation with the  $d \times d$  matrix valued kernel

$$(1.6) \quad K(x, y) = -(D^2 + \kappa^2)g_\kappa(x - y).$$

The integral operator thus defined involves second order distributional derivatives of the weakly singular kernel  $g_\kappa(x - y)$ . Instead of this form of an integro-differential operator, one can write the strongly singular integral operator also in the form of a Cauchy principal value integral, using the well-known relation (for more details, see section 3.5.1 below)

$$(1.7) \quad D^2 \int_{\mathbb{R}^d} g_\kappa(x - y)u(y)dy = \text{p.v.} \int_{\mathbb{R}^d} D^2 g_\kappa(x - y)u(y)dy - \frac{1}{d}u(x).$$

If we further assume that the permittivity  $\epsilon$  equals a constant  $\epsilon_r \in \mathbb{C} \setminus \{1\}$  in  $\Omega$ , we can divide by  $1 - \epsilon_r$  and arrive at the final form (1.1) with the integral understood in the principal value sense, the kernel given by (1.6), and the spectral parameter  $\lambda$  defined by the relation

$$(1.8) \quad \lambda = \frac{1}{1 - \epsilon_r} - \frac{1}{d} = \frac{d - 1 + \epsilon_r}{d(1 - \epsilon_r)}.$$

For  $d = 3$ , this relation  $\lambda = \frac{2 + \epsilon_r}{3(1 - \epsilon_r)}$  is known in the DDA literature as *Clausius-Mossotti polarizability*, referring to the fact that  $\frac{1}{\lambda}$  corresponds to the polarizability of the dipoles and to the Clausius-Mossotti equation between the molecular polarizability and the electric permittivity in a dielectric material, see for example [9, Section 4.5].

The principal part of the volume integral operator is obtained by taking the limit  $\kappa \rightarrow 0$ , and we will refer to this situation as the quasi-static Maxwell case. The resulting kernel is homogeneous of degree  $-d$ , and this property allows to analyze the corresponding linear system (1.2) using Fourier analysis of Toeplitz matrices. For this reason we study in this paper a class of strongly singular kernels that includes the quasi-static Maxwell kernel.

### 1.3. Outline of the paper.

In Section 1.5 we define a class of strongly singular kernels that are homogeneous of degree  $-d$  and translation invariant, and we evoke the relation between the numerical range of the corresponding singular integral operator in  $L^2$  and values of its symbol. The notion of numerical range allows to use a Lax-Milgram type argument to get a resolvent estimate for the restriction of the convolution operator to a bounded domain  $\Omega$ .

After introducing in Section 1.6 the delta-delta discretization, we state in Theorem 1.6 the main stability result valid for our class of operators.

In Section 2 we study tools for proving stability results, namely infinite Toeplitz matrices and their symbols defined by Fourier series. Here a main difficulty is that one needs precise bounds for the values of a function (numerical symbol) defined by a Fourier series that is not absolutely convergent. We find that one can use Ewald's method for this purpose. The result is that the symbol of the Toeplitz matrix is a bounded function, and that its range is always a superset of the range of the symbol of the integral operator, but that it might be strictly larger. If this is the case, then stability of the delta-delta scheme implies well-posedness of the integral equation, but not

vice versa: The numerical scheme does then not provide a spectrally correct approximation, and it might be unstable for values of the spectral parameter  $\lambda$  for which any Galerkin scheme of the integral equation would be stable.

In Section 3 we study in detail five representative examples.

Example 1 concerns the one-dimensional singular integral equation defined by the finite Hilbert transformation. Here the numerical symbol has a simple explicit expression, and this can be used to get estimates for the resolvent of the discretized operator by the resolvent of the integral operator, with constant equal to 1. This gives Theorem 3.4, which is the ideal stability result that subsequent results are measured against.

Examples 2 and 3 exhibit different behavior of the delta-delta scheme for two strongly singular integral operators in two dimensions. Whereas the two integral operators are equivalent, related by a simple rotation of the coordinate system, the two discrete systems show opposite behavior: We prove that in Example 2 the ranges of the symbol of the integral operator and of the symbol of the infinite Toeplitz matrix are identical, whereas in Example 3 there is an overshoot; the region of instability of the approximation scheme is strictly larger than the numerical range of the integral operator.

In Example 4, we graphically illustrate the relations, proved in Sections 1.5 and 2, between the spectrum and numerical range of the system matrices and the numerical range of the singular integral operator by considering a non-selfadjoint case. The kernel is a complex-valued function whose real and imaginary parts are given by the kernels of Examples 3 and 2, respectively.

The kernels studied in Examples 2 and 3 are also the off-diagonal and diagonal terms, respectively, in the matrix-valued kernel of the quasi-static Maxwell volume integral operator, which is the subject of Example 5. We study this for dimensions  $d \geq 2$  and give more precise results for  $d = 2$  and  $d = 3$ . In two dimensions we find the same overshoot of the numerical range of the numerical symbol versus the symbol of the integral operator as Example 3. In three dimensions this overshoot is even larger, and it can be verified numerically either by computing the numerical symbol using Ewald's method or by studying the asymptotic behavior of the smallest and largest eigenvalues of the matrix of the linear system as the mesh width tends to zero.

#### 1.4. Notation for Fourier transforms and Fourier series.

We use the following convention for the Fourier transformation in  $\mathbb{R}^d$ .

$$(1.9) \quad \widehat{f}(\xi) = \mathcal{F}f(\xi) = \int_{\mathbb{R}^d} f(x)e^{i\xi \cdot x} dx.$$

Inverse:

$$(1.10) \quad f(x) = \mathcal{F}^{-1}\widehat{f}(x) = (2\pi)^{-d} \int_{\mathbb{R}^d} \widehat{f}(\xi)e^{-ix \cdot \xi} d\xi.$$

For Fourier series, we use the following notation. For a sequence  $a : \mathbb{Z}^d \rightarrow \mathbb{C}$ , its Fourier series is defined as

$$(1.11) \quad \tilde{a}(\tau) = \sum_{m \in \mathbb{Z}^d} a(m)e^{im \cdot \tau}, \quad \tau \in Q = [-\pi, \pi]^d.$$

Inverse:

$$(1.12) \quad a(m) = (2\pi)^{-d} \int_Q \tilde{a}(\tau) e^{-im \cdot \tau} d\tau.$$

The definitions are extended in the usual way from convergent sums and integrals to suitable spaces of functions and distributions. In particular, we have Parseval's theorem

$$(1.13) \quad f \mapsto (2\pi)^{-\frac{d}{2}} \widehat{f} : L^2(\mathbb{R}^d) \rightarrow L^2(\mathbb{R}^d) \quad \text{and} \quad a \mapsto (2\pi)^{-\frac{d}{2}} \tilde{a} : \ell^2(\mathbb{Z}^d) \rightarrow L^2(Q)$$

are unitary (i.e. isometric Hilbert space isomorphisms).

Combining the Parseval formula and the convolution theorem gives

$$(1.14) \quad \int_{\mathbb{R}^{2d}} \overline{u(x)} k(x-y) v(y) dy dx = (2\pi)^{-d} \int_{\mathbb{R}^d} \overline{\widehat{u}(\xi)} \widehat{k}(\xi) \widehat{v}(\xi) d\xi,$$

$$(1.15) \quad \sum_{m,n \in \mathbb{Z}^d} \overline{a(m)} c(m-n) b(n) = (2\pi)^{-d} \int_Q \overline{\tilde{a}(\tau)} \tilde{c}(\tau) \tilde{b}(\tau) d\tau.$$

From these formulas follows immediately that the operators of convolution with  $k$  in  $L^2(\mathbb{R}^d)$  and of discrete convolution with  $c$  in  $\ell^2(\mathbb{Z}^d)$  are bounded if and only if the “symbols”  $\widehat{k}$  and  $\tilde{c}$  are bounded functions belonging to  $L^\infty(\mathbb{R}^d)$  and  $L^\infty(Q)$ , respectively.

Sufficient conditions for this are that  $k \in L^1(\mathbb{R}^d)$  and  $c \in \ell^1(\mathbb{Z}^d)$ . But these conditions are not necessary, and it is precisely the situation where they are not satisfied that will be relevant in the following.

We will use the *Poisson summation formula* in the form

$$(1.16) \quad \sum_{m \in \mathbb{Z}^d} f(m) e^{im \cdot \tau} = \sum_{n \in \mathbb{Z}^d} \widehat{f}(\tau + 2\pi n).$$

A sufficient (but in no way necessary) condition for (1.16) to hold for all  $\tau$  is that

$$f|_{\mathbb{Z}^d} \in \ell^1(\mathbb{Z}^d) \quad \text{and} \quad \widehat{f} \in L^1(\mathbb{R}^d).$$

If we do not assume  $f|_{\mathbb{Z}^d} \in \ell^1$ , but only  $\widehat{f} \in L^1$ , then  $f$  is bounded, the left hand side of (1.16) converges in the distributional sense and the right hand side converges in  $L^1(Q)$ . Then (1.16) is true in a weaker sense, the distributional left hand side being equal to the  $L^1(Q)$  right hand side.

*Example:* Gaussian with parameter  $s > 0$ .

$$(1.17) \quad f(x) = e^{-|x|^2 s} \quad \iff \quad \widehat{f}(\xi) = \left(\frac{\pi}{s}\right)^{\frac{d}{2}} e^{-\frac{|\xi|^2}{4s}}.$$

For this example the Poisson summation formula takes the form (for  $\tau \in \mathbb{R}^d$ )

$$(1.18) \quad \sum_{m \in \mathbb{Z}^d} e^{-|m|^2 s} e^{im \cdot \tau} = \sum_{n \in \mathbb{Z}^d} \left(\frac{\pi}{s}\right)^{\frac{d}{2}} e^{-\frac{|\tau + 2\pi n|^2}{4s}}.$$

## 1.5. Kernels and their symbols.

1.5.1. *Homogeneous kernels.* Later on, we will consider a rather restricted class of strongly singular integral operators on  $\mathbb{R}^d$  that are convolutions with kernel functions of the form

$$(1.19) \quad K(x) = p(x) |x|^{-d-2} \quad \text{where } p \text{ is a homogeneous polynomial of degree } 2.$$

But first we recall some well-known general properties of homogeneous functions and distributions that can be found, for example, in [8, Chap. III].

Let  $K$  be a function on  $\mathbb{R}^d$ , positively homogeneous of degree  $-d$  and smooth outside of the origin. For a given  $\epsilon > 0$ , one can define a distribution  $K_\epsilon \in \mathcal{S}'(\mathbb{R}^d)$  that coincides with  $K$  on  $\mathbb{R}^d \setminus \{0\}$  by its action on a test function  $\phi$  as

$$(1.20) \quad \langle K_\epsilon, \phi \rangle = \int_{|x| < \epsilon} K(x)(\phi(x) - \phi(0)) dx + \int_{|x| > \epsilon} K(x)\phi(x) dx.$$

This is independent of  $\epsilon$  if and only if  $K$  satisfies the cancellation condition on the unit sphere  $\mathbb{S}^{d-1}$

$$(1.21) \quad \int_{\mathbb{S}^{d-1}} K ds = 0.$$

In this case, we denote the distribution simply by  $K$ , and we can take the limit  $\epsilon \rightarrow 0$ , thus we get the Cauchy principal value.

$$(1.22) \quad \langle K, \phi \rangle = \text{p.v.} \int K(x)\phi(x) dx = \lim_{\epsilon \rightarrow 0} \int_{|x| > \epsilon} K(x)\phi(x) dx.$$

Another consequence of the cancellation condition (1.21) is that the Fourier transform  $\widehat{K}$  of the homogeneous distribution  $K$  is a bounded function homogeneous of degree 0, smooth outside of the origin and also satisfying the cancellation condition. The operator  $A$  of convolution with  $K$  is therefore bounded in  $L^2(\mathbb{R}^d)$ . Note that in the absence of condition (1.21),  $\widehat{K}_\epsilon$  would have a logarithmic singularity at 0.

The operator  $A$  is diagonalized by Fourier transformation:

$$(1.23) \quad \mathcal{F} Au = \widehat{K} \widehat{u}.$$

Therefore in  $L^2(\mathbb{R}^d)$ , we can obtain information about the spectrum  $\text{Sp}(A)$  and about the numerical range  $W(A)$  from the corresponding easily checked information about the operator of multiplication by the symbol  $\widehat{K}$ .

We recall that the numerical range of  $A$  is defined by

$$W(A) = \{(u, Au) \mid \|u\| = 1\},$$

where  $(\cdot, \cdot)$  denotes the Hilbert space inner product. It is convex by the Toeplitz-Hausdorff theorem and it contains the spectrum of  $A$ . Denote by  $\text{im}(\widehat{K}) = \widehat{K}(\mathbb{R}^d)$  the image (range) of  $\widehat{K}$ . This is a compact set. We note a first result implied by the unitary equivalence (1.23) with the multiplication operator.

**Lemma 1.1.** *The spectrum  $\text{Sp}(A)$  is the image  $\text{im}(\widehat{K})$ , and the closure  $\overline{W(A)}$  of the numerical range of  $A$  is the convex hull of  $\text{im}(\widehat{K})$ .*



It is well known (and easy to prove) that the numerical range allows estimates for the operator norm of the resolvent: For any  $\lambda \in \mathbb{C} \setminus \overline{W(A)}$ ,

$$\|(\lambda\mathbb{I} - A)^{-1}\| \leq \text{dist}(\lambda, W(A))^{-1}.$$

It is also monotone with respect to inclusions of subspaces, a property not shared by the spectrum. Given an open set  $\Omega \subset \mathbb{R}^d$ , we denote by  $A_\Omega$  the restriction of the convolution operator  $A$  to  $L^2(\Omega)$  and consider the strongly singular integral equation  $(\lambda\mathbb{I} - A_\Omega)u = f$ , or in detail

$$(1.24) \quad \lambda u(x) - \text{p.v.} \int_{\Omega} K(x-y)u(y) dy = f(x) \quad (x \in \Omega).$$

From the definition of the numerical range follows immediately the inclusion  $W(A_\Omega) \subset W(A)$ .

We can summarize this discussion:

**Proposition 1.2.** *Let  $\mathcal{C} \subset \mathbb{C}$  be a closed convex set such that  $\widehat{K}(\xi) \in \mathcal{C}$  for all  $\xi \in \mathbb{S}^{d-1}$ . Then for all  $\lambda \notin \mathcal{C}$  and any  $f \in L^2(\Omega)$ , the integral equation (1.24) has a unique solution  $u \in L^2(\Omega)$ , and there is a resolvent estimate in the  $L^2(\Omega)$  norm*

$$(1.25) \quad \|(\lambda\mathbb{I} - A_\Omega)^{-1}\| \leq \text{dist}(\lambda, \mathcal{C})^{-1}.$$

*Remark 1.3.* The same argument implies stability for any Galerkin method: Let  $X_h$  be any closed subspace of  $L^2(\Omega)$ , and let  $A_h : X_h \rightarrow X_h$  be the operator defined by restricting the sesquilinear form  $(u, Av)$  to  $X_h \times X_h$ . Then the statement of Proposition 1.2 remains true if we replace  $A_\Omega$  by  $A_h$ .

*Remark 1.4.* Whereas there is, in general, no simple relation between the spectra  $\text{Sp}(A_\Omega)$  and  $\text{Sp}(A)$ , for the numerical ranges of our convolution operators with homogeneous kernels we not only have the inclusion  $W(A_\Omega) \subset W(A)$ , but also the converse. Namely there holds

$$(1.26) \quad \overline{W(A_\Omega)} = \overline{W(A)} \quad \text{for any open subset } \Omega \subset \mathbb{C}.$$

*Proof.* The set of Rayleigh quotients  $\frac{(u, Au)}{(u, u)}$ , where  $u \in L^2(\mathbb{R}^d) \setminus \{0\}$  has compact support, is a dense subset of  $W(A)$ . We show that it is a subset of  $W(A_\Omega)$ : Indeed, let  $u$  be such a function and let  $\rho > 0$  and  $a \in \mathbb{R}^d$  be chosen such that the support of the function  $u_{\rho, a}$  defined by  $u_{\rho, a}(x) = u(\rho x + a)$  is contained in  $\Omega$ . Then

$$\frac{(u, Au)}{(u, u)} = \frac{(u_{\rho, a}, Au_{\rho, a})}{(u_{\rho, a}, u_{\rho, a})} \in W(A_\Omega).$$

□

1.5.2. *Special kernels.* For  $d = 1$ , there is essentially only one non-trivial kernel homogeneous of degree  $-d$ , namely  $K(x) = \frac{1}{x}$ .

In  $\mathbb{R}^d$  for  $d \geq 2$ , while some of the following analysis would be possible for more general homogeneous kernels, we focus now on the situation (1.19). This means that from now on, we fix a strongly singular kernel  $K$  and a homogeneous polynomial  $p$  of degree 2 with  $K(x) = p(x)|x|^{-d-2}$ , satisfying (1.21), considered as a distribution on  $\mathbb{R}^d$  according to (1.22), and we denote by  $\widehat{K}$  its Fourier transform.

**Lemma 1.5.** *Let  $K$  have the form (1.19) and satisfy (1.21). Then*

$$(1.27) \quad \widehat{K}(\xi) = -\nu_d \frac{p(\xi)}{|\xi|^2}, \quad \text{where } \nu_d = \frac{2\pi^{\frac{d}{2}}}{d\Gamma(\frac{d}{2})} \text{ is the volume of the unit ball in } \mathbb{R}^d.$$

*Proof.* We first compute the Fourier transform of  $p(x)e^{-|x|^2s}$ , using (1.17)

$$\mathcal{F}_{x \rightarrow \xi}[p(x)e^{-|x|^2s}] = \left(\frac{\pi}{s}\right)^{\frac{d}{2}} p(-i\partial_\xi) e^{-\frac{|\xi|^2}{4s}}.$$

The evaluation of these derivatives leads to the following simple result, as we will show:

$$(1.28) \quad \mathcal{F}_{x \rightarrow \xi}[p(x)e^{-|x|^2s}] = -\left(\frac{\pi}{s}\right)^{\frac{d}{2}} \frac{1}{4s^2} p(\xi) e^{-\frac{|\xi|^2}{4s}}.$$

For  $j, k \in \{1, \dots, d\}$  with  $j \neq k$ , let

$$(1.29) \quad a_{jk}(x) = x_j^2 - x_k^2, \quad b_{jk}(x) = x_j x_k.$$

Any homogeneous polynomial of degree 2 satisfying the cancellation condition  $\int_{\mathbb{S}^{d-1}} p = 0$  is a linear combination of the  $a_{jk}$  and  $b_{jk}$ , so we need to verify (1.28) only for these.

Note that  $\partial_{\xi_j} e^{-\frac{|\xi|^2}{4s}} = -\frac{1}{2s} \xi_j e^{-\frac{|\xi|^2}{4s}}$  and  $\partial_{\xi_j}^2 e^{-\frac{|\xi|^2}{4s}} = \left(-\frac{1}{2s} + \frac{1}{4s^2} \xi_j^2\right) e^{-\frac{|\xi|^2}{4s}}$ .

Then for  $p = a_{jk}$ , we see

$$(\partial_{\xi_k}^2 - \partial_{\xi_j}^2) e^{-\frac{|\xi|^2}{4s}} = \frac{1}{4s^2} (\xi_k^2 - \xi_j^2) e^{-\frac{|\xi|^2}{4s}},$$

and for  $p = b_{jk}$ , we see

$$\partial_{\xi_j} \partial_{\xi_k} e^{-\frac{|\xi|^2}{4s}} = \frac{1}{4s^2} \xi_k \xi_j e^{-\frac{|\xi|^2}{4s}}.$$

Thus in both cases we have

$$(1.30) \quad p(\partial_\xi) e^{-\frac{|\xi|^2}{4s}} = \frac{1}{4s^2} p(\xi) e^{-\frac{|\xi|^2}{4s}},$$

and (1.28) is proved.

Now we use the definition of the Gamma function

$$\Gamma(a) = \int_0^\infty t^a e^{-t} \frac{dt}{t} = |x|^{2a} \int_0^\infty s^a e^{-|x|^2s} \frac{ds}{s}$$

to write the kernel as an integral over Gaussians

$$(1.31) \quad K(x) = p(x) |x|^{-d-2} = \frac{1}{\Gamma(\frac{d}{2} + 1)} \int_0^\infty s^{\frac{d}{2}+1} p(x) e^{-|x|^2s} \frac{ds}{s}.$$

Taking Fourier transforms and using (1.28), we find with  $u = \frac{|\xi|^2}{4s}$

$$(1.32) \quad -\widehat{K}(\xi) = \frac{\pi^{\frac{d}{2}}}{4\Gamma(\frac{d}{2} + 1)} \int_0^\infty s^{-1} p(\xi) e^{-\frac{|\xi|^2}{4s}} \frac{ds}{s} = p(\xi) |\xi|^{-2} \frac{\pi^{\frac{d}{2}}}{\frac{d}{2}\Gamma(\frac{d}{2})} \int_0^\infty e^{-u} du = \nu_d p(\xi) |\xi|^{-2}$$

as claimed.  $\square$

### 1.6. Delta-delta discretization.

Let  $N \in \mathbb{N}$ , fix some origin  $a^N \in \mathbb{R}^d$  and define the cubic grid of meshwidth  $h = \frac{1}{N}$  by

$$\Sigma^N = \{x_m^N = a^N + \frac{m}{N} \mid m \in \mathbb{Z}^d\}.$$

We further define

$$\omega^N = \{m \in \mathbb{Z}^d \mid x_m^N \in \Omega\}.$$

Then a very simple discretization of the strongly singular integral equation  $(\lambda\mathbb{I} - A_\Omega)u = f$  (1.24) is the following

$$(1.33) \quad \lambda u_m - N^{-d} \sum_{n \in \omega^N, m \neq n} K(x_m^N - x_n^N) u_n = f(x_m^N), \quad (m \in \omega^N),$$

or in shorthand  $(\lambda\mathbb{I} - T^N)U = F$ .

The name ‘‘delta-delta discretization’’ points at the fact that this discretization formally looks like a Galerkin method for the integral equation (1.24) with Dirac deltas as both test and trial functions, except for the diagonal terms of  $T^N$ , where we put zero, which is natural in view of the cancellation condition (1.21).

Our aim in this paper is to analyze the linear system (1.33), in particular its stability in  $\ell^2(\omega^N)$ , in the same way as we did above for the integral equation (1.24) in  $L^2(\Omega)$ , and to compare the two.

We state a general result here, which we prove in the next section. More precise results will be given below in Section 3 for some examples, in particular those mentioned in Subsection 1.1.

**Theorem 1.6.** *Let  $K$  be a strongly singular kernel satisfying (1.19) and (1.21). Then there exists a compact convex set  $\mathcal{C} \subset \mathbb{C}$  such that for any  $\lambda \in \mathbb{C} \setminus \mathcal{C}$ , any  $N \in \mathbb{N}$  for which  $\omega^N$  is non-empty, and for any  $F \in \ell^2(\Omega^N)$ , the system (1.33) has a unique solution, and there is a uniform estimate for the inverse in the  $\ell^2(\Omega^N)$  operator norm*

$$(1.34) \quad \|(\lambda\mathbb{I} - T^N)^{-1}\|_{\mathcal{L}(\ell^2(\Omega^N))} \leq \text{dist}(\lambda, \mathcal{C})^{-1}.$$

Furthermore, with the strongly singular integral operator  $A$  defined above in Section 1.5.1, there holds the inclusion

$$(1.35) \quad W(A) \subset \mathcal{C}.$$

*Remark 1.7.* Note that the inclusion  $W(A) \subset \mathcal{C}$  implies that for  $\lambda \notin \mathcal{C}$  the singular integral equation is uniquely solvable, too, and provides the a priori estimate (1.25), for any domain  $\Omega \subset \mathbb{R}^d$ . On the other hand, in order to guarantee stability for  $\lambda \in \mathbb{C} \setminus \mathcal{C}$ , the inclusion may need to be strict, as we shall see in the examples, and then there may be  $\lambda \in \mathcal{C} \setminus W(A)$  for which the singular integral equation is well posed, but the delta-delta discretization is *unstable*.

## 2. THE DISCRETE SYSTEM

Let  $T^N$  be the matrix representing the discretized integral operator in (1.33):

$$(2.1) \quad T^N = (t_{mn}^N)_{m,n \in \omega^N} \quad \text{with} \quad t_{mn}^N = \begin{cases} N^{-d} K(x_m^N - x_n^N) & (m \neq n) \\ 0 & (m = n) \end{cases}.$$

Our aim is to bound the numerical range  $W(T^N)$  independently of  $N$ .

### 2.1. Toeplitz structure.

The matrix elements  $t_{mn}^N$  of  $T^N$  do not depend on the choice of the origin  $a^N$ , and since we assumed that  $K$  is homogeneous of degree  $-d$ , we have

$$N^{-d} K(x_m^N - x_n^N) = K(m - n),$$

hence  $T^N$  is a finite section of a fixed infinite Toeplitz (discrete convolution) matrix

$$(2.2) \quad T = (t_{mn})_{m,n \in \mathbb{Z}^d} \quad \text{with} \quad t_{mn} = \begin{cases} K(m - n) & (m \neq n), \\ 0 & (m = n). \end{cases}$$

Theorem 1.6 will be proved if we can show that  $T$  defines a bounded linear operator in  $\ell^2(\mathbb{Z}^d)$  whose numerical range  $W(T)$  contains  $W(A)$ . We can then choose  $\mathcal{C}$  as the closure of  $W(T)$ .

We use Fourier series and the convolution theorem to diagonalize the matrix  $T$  and to represent the sesquilinear form defined by the matrix  $T^N$ , compare (1.15). For  $U = (u_m)_{m \in \omega^N}$ , we find

$$(2.3) \quad (U, T^N U) = (2\pi)^{-d} \int_Q F(\tau) |\tilde{u}(\tau)|^2 d\tau.$$

Here  $\tilde{u}(\tau) = \sum_{m \in \omega^N} u_m e^{im \cdot \tau}$  and  $Q = [-\pi, \pi]^d$ .  $F(\tau)$  is the symbol (characteristic function) of the Toeplitz matrix  $T$ :

$$(2.4) \quad F(\tau) = \sum_{m \in \mathbb{Z}^d, m \neq 0} K(m) e^{im \cdot \tau}.$$

The problem is now reduced to the study of the operator of multiplication by the function  $F$  in  $L^2(Q)$ .

**Lemma 2.1.** *The operator  $T : \ell^2(\mathbb{Z}^d) \rightarrow \ell^2(\mathbb{Z}^d)$  is bounded if and only if  $F \in L^\infty(Q)$ .*

*The closure of  $W(T)$  is the closed convex hull of the range  $\text{im}(F) = \{F(\tau) \mid \tau \in Q\}$  and is also equal to the closure of the union  $\bigcup_{N \in \mathbb{N}} W(T^N)$ .*

Proof: This is immediate from (2.3).

### 2.2. Ewald method.

The problem that makes the statement  $F \in L^\infty(Q)$  non trivial is that the Fourier series (2.4) is not absolutely convergent. The sequence  $(K(m))_{m \in \mathbb{Z}^d}$  is of order  $O(|m|^{-d})$  at infinity and therefore in  $\ell^p(\mathbb{Z}^d)$  for all  $p > 1$ , but not for  $p = 1$ . Its membership in  $\ell^2(\mathbb{Z}^d)$  implies, for example, that the series converges in the sense of  $L^2(Q)$ . The slow convergence of the Fourier series for  $F$  makes it also unsuitable for using it in numerical computations to find bounds for  $\text{im}(F)$ .

We will use a variant of a method introduced by P. P. Ewald [7] in 1921 as a tool to compute slowly converging lattice sums. It has become a routine method for the computation of periodic and quasi-periodic Green functions, with application in numerical electrodynamics and other fields where periodic structures appear. Among the many presentations of the method: Appendix A of the article [6] or Section 2.13.3 in the book [2].

We use it here as a summation method for our slowly converging Fourier series. In our restricted setting it turns out to give surprisingly simple results.

The method introduces a decomposition  $K = K^F + K^P$  for the coefficients and correspondingly  $F = F^F + F^P$  for the Fourier series in such a way that both  $K^F$  and the Fourier transform  $\widehat{K}^P$  of  $K^P$  are exponentially decreasing at infinity, so that both the Fourier series for  $F^F(\tau)$  and the Poisson sum (compare (1.16)) for  $F^P(\tau)$  are rapidly convergent, which not only proves the boundedness of  $F$ , but gives also a fast numerical algorithm for its computation.

In the literature one often labels the two terms in the decomposition “spatial” and “spectral” sums, but this is not pertinent to our situation, where the lattice sum runs over the Fourier variable, and the Fourier series runs over spatial points. So we will use “Fourier” and “Poisson” sums as labels.

The idea of Ewald’s method is to represent  $K(x)$  by an integral over Gaussians from 0 to  $\infty$  as we did already in Section 1.5.2 above:

$$(2.5) \quad K(x) = p(x)|x|^{-d-2} = \frac{p(x)}{\Gamma(\frac{d}{2} + 1)} \int_0^\infty s^{\frac{d}{2}} e^{-|x|^2 s} ds$$

and then to split the integral at a point  $\beta^2 > 0$ :

$$(2.6) \quad K^F(x) = \frac{p(x)}{\Gamma(\frac{d}{2} + 1)} \int_{\beta^2}^\infty s^{\frac{d}{2}} e^{-|x|^2 s} ds,$$

$$(2.7) \quad K^P(x) = \frac{p(x)}{\Gamma(\frac{d}{2} + 1)} \int_0^{\beta^2} s^{\frac{d}{2}} e^{-|x|^2 s} ds.$$

We see that  $K^F$  is simply the product of  $K$  by a function exponentially decreasing at infinity

$$(2.8) \quad K^F(x) = K(x) \frac{\Gamma(\frac{d}{2} + 1, \beta^2|x|^2)}{\Gamma(\frac{d}{2} + 1)}$$

with the (upper) incomplete Gamma function (see [1, §6.5])

$$\Gamma(a, x) = \int_x^\infty t^{a-1} e^{-t} dt.$$

Therefore  $K^F(x) = O(|x|^2 e^{-\beta^2|x|^2})$  as  $|x| \rightarrow \infty$ , and the Fourier series for  $F^F(\tau)$

$$(2.9) \quad F^F(\tau) = \sum_{m \in \mathbb{Z}^d, m \neq 0} K^F(m) e^{im \cdot \tau}$$

converges rapidly, implying that  $F^F$  is an analytic function on  $\mathbb{R}^d / (2\pi\mathbb{Z})^d$ .

Consequently, the Fourier series for  $F^P(\tau)$  converges as slowly as the one for  $F(\tau)$ , and we use instead the Poisson summation formula (1.16) and write

$$(2.10) \quad F^P(\tau) = \sum_{n \in \mathbb{Z}^d} \widehat{K}^P(\tau + 2\pi n).$$

We can evaluate  $\widehat{K}^P$  with the formulas used for  $\widehat{K}$  in Lemma 1.5. As in (1.32) we obtain

$$(2.11) \quad \begin{aligned} \widehat{K}^P(\xi) &= \frac{-\pi^{\frac{d}{2}}}{4\Gamma(\frac{d}{2}+1)} \int_0^{\beta^2} s^{-1} p(\xi) e^{-\frac{|\xi|^2}{4s}} \frac{ds}{s} = -p(\xi) |\xi|^{-2} \frac{\pi^{\frac{d}{2}}}{2\Gamma(\frac{d}{2})} \int_{\frac{|\xi|^2}{4\beta^2}}^{\infty} e^{-u} du \\ &= \widehat{K}(\xi) e^{-\frac{|\xi|^2}{4\beta^2}}. \end{aligned}$$

Therefore we also obtain a very simple form for the Fourier transform, namely that  $\widehat{K}^P$  is just the symbol of  $A$  cut off at infinity, and therefore the series (2.10) converges absolutely and uniformly. At most one term in the sum may be discontinuous, when  $\tau + 2\pi n = 0$ , and for  $\tau \in Q$  this is the term with  $n = 0$ . We can summarize the result.

**Proposition 2.2.** *The symbol  $F(\tau)$  of the infinite Toeplitz matrix  $T$  is a bounded function given for any  $\beta > 0$  by the exponentially convergent sums*

$$(2.12) \quad F(\tau) = \sum_{m \in \mathbb{Z}^d, m \neq 0} K(m) \frac{\Gamma(\frac{d}{2}+1, \beta^2 |m|^2)}{\Gamma(\frac{d}{2}+1)} e^{im \cdot \tau} + \sum_{n \in \mathbb{Z}^d} \widehat{K}(\tau + 2\pi n) e^{-\frac{|\tau+2\pi n|^2}{4\beta^2}}.$$

In the period cube  $Q = [-\pi, \pi]^d$ , it is  $C^\infty$  outside of 0, and it has the form

$$(2.13) \quad F(\tau) = \widehat{K}(\tau) + F_0(\tau) \quad \text{where } F_0 \text{ is analytic in } Q \text{ and } F_0(0) = 0.$$

*Proof.* We have proved equation (2.12) above, except for one point: From Poisson's summation formula follows that the Poisson sum (2.10) equals the Fourier series with coefficients  $K^P(m)$ ,  $m \in \mathbb{Z}^d$ , including  $m = 0$ . But in the Fourier series (2.4) defining  $F(t)$  as well as in (2.9) defining  $F^F(t)$ , we have excluded  $m = 0$ . So we should compensate for  $K^P(0)$ , which according to (2.7) equals

$$K^P(0) = \frac{p(0)\beta^d}{\Gamma(\frac{d}{2}+1)}.$$

Now, since we assumed  $p(x)$  to be a homogeneous polynomial of degree 2, we have  $p(0) = 0$  and hence no compensation is needed.

Representing  $F_0$  as

$$F_0(\tau) = \sum_{m \in \mathbb{Z}^d, m \neq 0} K(m) \frac{\Gamma(\frac{d}{2}+1, \beta^2 |m|^2)}{\Gamma(\frac{d}{2}+1)} e^{im \cdot \tau} + \sum_{n \in \mathbb{Z}^d, n \neq 0} \widehat{K}(\tau + 2\pi n) e^{-\frac{|\tau+2\pi n|^2}{4\beta^2}} + \widehat{K}(\tau) (e^{-\frac{|\tau|^2}{4\beta^2}} - 1),$$

we see immediately that it is analytic. For finding  $F_0(0)$ , we can use the following observation.

**Lemma 2.3.** *Let  $S \subset \mathbb{R}^d$  be a finite set that is cubically symmetric, i. e. invariant under reflections at coordinate planes and under permutations of the coordinates, and let  $p$  be a homogeneous polynomial of degree 2 satisfying the cancellation condition  $\int_{\mathbb{S}^{d-1}} p = 0$ . Then*

$$\sum_{x \in S} p(x) = 0.$$

This is immediately clear when  $p$  is one of the  $a_{jk}$  or  $b_{jk}$  from (1.29), and it is therefore true for all  $p$  satisfying the (spherical) cancellation condition.

For any  $M \in \mathbb{R}$ , the set  $\{m \in \mathbb{Z}^d \mid |m|^2 = M\}$  is either empty or cubically symmetric. Therefore for  $\tau = 0$ , the two sums in the representation of  $F_0(\tau)$  are 0. The last term

$$\widehat{K}(\tau) \left( e^{-\frac{|\tau|^2}{4\beta^2}} - 1 \right) = -\nu_d p(\tau) \frac{e^{-\frac{|\tau|^2}{4\beta^2}} - 1}{|\tau|^2}$$

tends to  $\nu_D p(0)/(4\beta^2) = 0$  as  $\tau \rightarrow 0$ , and hence  $F_0(0) = 0$ .  $\square$

Proposition 2.2 implies Theorem 1.6, where  $\mathcal{C}$  is the closed convex hull of  $\text{im}(F)$ . The inclusion  $W(A) \subset \mathcal{C}$  is easy to see from (2.13):

Given  $\epsilon > 0$ , let  $\delta > 0$  be such that for  $|\tau| < \delta$  we have  $|F_0(\tau)| < \epsilon$ . Since  $\widehat{K}$  is homogeneous of degree zero, it takes all of its values already on the ball  $B_\delta(0)$  of radius  $\delta$ . Thus

$$\text{im}(\widehat{K}) \subset F(B_\delta(0)) + B_\epsilon(0) \subset \text{im}(F) + B_\epsilon(0).$$

Taking convex hulls shows that

$$W(A) \subset \mathcal{C} + B_\epsilon(0) \quad \text{for all } \epsilon > 0.$$

*Remark 2.4.* The very simple form of the Ewald representation (2.12) comes from the very simple form of the Fourier transforms (1.27) and (1.28), which in turn rely on the cancellation condition (1.21). Now for the kernel  $K$  this condition is natural, because it is necessary in order to represent  $K$  as a homogeneous distribution and to have a bounded Fourier transform. But for the symbol  $\widehat{K}$  it is not as natural. We can add a constant and still have a function homogeneous of degree zero, which will then not satisfy the cancellation condition. An example is  $\xi_j \xi_k |\xi|^{-2}$  for all  $j, k$ , even for  $j = k$ .

On the other hand, the representation  $K(x) = p(x)|x|^{-d-2}$  may not be the most natural, one may come across cases (see Example 5 below) like

$$K_{jk}(x) = \delta_{jk}|x|^{-d} - d x_j x_k |x|^{-d-2},$$

where for  $j = k$  the two terms in the sum do not separately satisfy (1.21). This fits into our framework, however, because

$$K_{jk}(x) = -d b_{jk}(x) |x|^{-d-2} \quad \text{for } j \neq k, \quad \text{and} \quad K_{kk}(x) = \sum_{j=1}^d a_{jk}(x) |x|^{-d-2}.$$

If one treats the two terms individually, one may get formulas for Fourier transforms and for the Ewald splitting that are less symmetric than what we presented above.

### 2.3. An integral representation.

We have another look at the Ewald splitting for the numerical symbol  $F(\xi) = F^F(\xi) + F^P(\xi)$  described in (2.6)–(2.12)

$$(2.14) \quad F^F(\xi) = \sum_{m \in \mathbb{Z}^d} \frac{p(m)}{\Gamma(\frac{d}{2} + 1)} \int_{\beta^2}^{\infty} s^{\frac{d}{2}} e^{-|m|^2 s} ds e^{im \cdot \xi}$$

$$(2.15) \quad F^P(\xi) = \sum_{n \in \mathbb{Z}^d} \frac{-\pi^{\frac{d}{2}}}{4\Gamma(\frac{d}{2} + 1)} \int_0^{\beta^2} s^{-2} p(\xi + 2\pi n) e^{-\frac{|\xi + 2\pi n|^2}{4s}} ds.$$

These formulas are valid for any  $0 < \beta < \infty$ . All the sums and integrals are converging absolutely here, and therefore we can interchange sums and integrals.

$$(2.16) \quad F^F(\xi) = \int_{\beta^2}^{\infty} H^F(\xi, s) ds \quad \text{with} \quad H^F(\xi, s) = \sum_{m \in \mathbb{Z}^d} \frac{p(m)}{\Gamma(\frac{d}{2} + 1)} s^{\frac{d}{2}} e^{-|m|^2 s} e^{im \cdot \xi}$$

$$(2.17) \quad F^P(\xi) = \int_0^{\beta^2} H^P(\xi, s) ds \quad \text{with} \quad H^P(\xi, s) = \sum_{n \in \mathbb{Z}^d} \frac{-\pi^{\frac{d}{2}} p(\xi + 2\pi n)}{4\Gamma(\frac{d}{2} + 1)} s^{-2} e^{-\frac{|\xi + 2\pi n|^2}{4s}}.$$

From the definition (2.16) of  $H^F$  and the fact that  $|m| \geq 1$  in the sum follows without difficulty that for any  $0 < \gamma < 1$  there exists a constant  $C$  such that

$$(2.18) \quad |H^F(\xi, s)| \leq C e^{-\gamma s} \quad \text{for all } s \geq 1, \xi \in \mathbb{R}^d.$$

To see the behavior of  $H^P(\xi, s)$  from (2.17), we decompose

$$H^P(\xi, s) = H_0(\xi, s) + H_1(\xi, s)$$

with

$$(2.19) \quad H_0(\xi, s) = -\frac{\pi^{\frac{d}{2}}}{4\Gamma(\frac{d}{2} + 1)} \sum_{n \in \mathbb{Z}^d, n \neq 0} p(\xi + 2\pi n) s^{-2} e^{-\frac{|\xi + 2\pi n|^2}{4s}},$$

$$(2.20) \quad H_1(\xi, s) = -\frac{\pi^{\frac{d}{2}}}{4\Gamma(\frac{d}{2} + 1)} p(\xi) s^{-2} e^{-\frac{|\xi|^2}{4s}}.$$

Now we use the fact that for  $\xi \in Q$  and  $n \neq 0$  we have  $|\xi + 2\pi n| \geq \pi$ . Therefore for any  $\delta < \frac{\pi^2}{4}$  there is a constant  $C$  such that

$$(2.21) \quad |H_0(\xi, s)| \leq C e^{-\frac{\delta}{s}} \quad \text{for all } 0 < s \leq 1, \xi \in Q,$$

and  $H_0(\xi, s)$  is analytic in  $\xi$  for all  $s$ .

It remains to analyze the term with  $n = 0$ , i.e.  $H_1$ . It is clear that it vanishes for  $\xi = 0$ , and for every  $\xi \neq 0$  there exists a constant  $C_\xi$  and  $0 < \gamma < \frac{|\xi|^2}{4}$  such that

$$(2.22) \quad |H_1(\xi, s)| \leq C_\xi \min\{s^{-2}, e^{-\frac{\gamma}{s}}\} \quad \text{for all } s \in (0, \infty).$$

Thus  $H_1(\xi, s)$  is integrable over  $s \in (0, \infty)$  for all  $\xi$ , but there is no uniform bound for  $C_\xi$ : Considering  $\sup_{s>0} s^{-2} e^{-\frac{|\xi|^2}{4s}}$ , one sees that  $C_\xi = O(|\xi|^{-2})$  as  $\xi \rightarrow 0$ .

Thus we see that  $H^F$  is integrable as  $s \rightarrow \infty$  according to (2.18), and  $H^P$  is integrable as  $s \rightarrow 0$  according to (2.21) and (2.22), but, because of Poisson's summation formula, they are in fact the same

$$H^F(\xi, s) = H^P(\xi, s),$$

so we can use all of the above estimates for both of them. We can summarize

**Proposition 2.5.** *The symbol  $F(\xi)$  has the integral representation*

$$(2.23) \quad F(\xi) = \int_0^\infty H(\xi, s) ds,$$



where  $H(\xi, s)$  is given either by the Fourier series  $H^F$  in (2.16) or, equivalently, by the lattice sum  $H^P$  in (2.17). The decomposition  $F = F_0 + \widehat{K}$  in Proposition 2.2 corresponds to the decomposition  $H = H_0 + H_1$  with  $H_0$  and  $H_1$  defined in (2.19) and (2.20), and there holds

$$(2.24) \quad F_0(\xi) = \int_0^\infty H_0(\xi, s) ds \quad \text{and} \quad \widehat{K}(\xi) = \int_0^\infty H_1(\xi, s) ds.$$

In these integrals, the functions  $s \mapsto H_0(\xi, s)$ ,  $s \mapsto H_1(\xi, s)$ , and  $s \mapsto H(\xi, s)$  are integrable on  $(0, \infty)$  for any  $\xi \in Q$ , for any  $\xi \in \mathbb{R}^d \setminus \{0\}$ , and for any  $\xi \in Q \setminus \{0\}$ , respectively.

The integral representations (2.23) and (2.24) will be used below to get bounds for the function  $F(\xi)$  from estimates for  $H(\xi, s)$ . The latter will be a consequence of the following observation that can be proved using Fourier representations (2.16) for  $H$  and (1.28) for  $H_1$ .

**Lemma 2.6.** *The functions*

$$(\xi, s) \mapsto s^{-\frac{d}{2}} H_0(\xi, s), \quad (\xi, s) \mapsto s^{-\frac{d}{2}} H_1(\xi, s), \quad (\xi, s) \mapsto s^{-\frac{d}{2}} H(\xi, s),$$

are solutions of the heat equation

$$(\partial_s - \Delta_\xi)u(\xi, s) = 0 \quad \text{in } Q \times (0, \infty).$$

#### 2.4. Matrix-valued kernels.

Until now, we have considered kernel functions with values in  $\mathbb{C}$  and integral operators acting on scalar functions. The generalization to vector-valued functions and matrix-valued kernels is simple and straightforward, and we do not find it necessary to introduce typographic distinctions for the vector-valued objects. The main difference is that in the general theory of Section 1.5, one has to use the numerical range  $W(K(x))$  of the matrix  $K(x)$  instead of the value  $K(x)$  in statements such as Lemma 1.1 and Proposition 1.2. In particular

$$(2.25) \quad \overline{W(A)} \text{ is the closed convex hull of } \bigcup_{\xi \in \mathbb{R}^d} W(\widehat{K}(\xi)).$$

Theorem 1.6 remains literally true, but for the construction of the set  $\mathcal{C}$  one has once again to use the numerical range  $W(F(t))$  of the matrix-valued function  $F$ . In Lemma 2.1, the characterization of the numerical range  $W(T)$  is to be understood as follows.

**Lemma 2.7.** *The closure of  $W(T)$  is the closed convex hull of  $\bigcup_{\tau \in Q} W(F(\tau))$  and is also equal to the closure of the union  $\bigcup_{N \in \mathbb{N}} W(T^N)$ .*

The basic Parseval-convolution formula (1.15) now has to be written, instead of the scalar version (2.3), as

$$(2.26) \quad (U, T^N U) = (U, T U) = (2\pi)^{-d} \int_Q \overline{\tilde{u}(\tau)}^\top F(\tau) \tilde{u}(\tau) d\tau.$$

Here  $\tilde{u}(\tau) = \sum_{m \in \omega^N} u_m e^{im \cdot \tau}$ , and  $F(\tau)$  is the matrix-valued symbol of the block Toeplitz matrix  $T = (K(m-n))_{m, n \in \mathbb{Z}^d}$ :

$$(2.27) \quad F(\tau) = \sum_{m \in \mathbb{Z}^d, m \neq 0} K(m) e^{im \cdot \tau}.$$

From (2.26) one can immediately read the properties of the numerical range stated in Lemma 2.7.

In this paper, most considered examples of kernels are real-valued and the matrices symmetric, in which case the integral operators are selfadjoint, and the numerical ranges consist of intervals in the real line.

### 3. EXAMPLES

#### 3.1. Example 1. Dimension $d = 1$ . Finite Hilbert transformation.

We start with the simplest example of a strongly singular integral equation and show that the stability of its delta-delta approximation can be completely analyzed, resulting in a kind of ideal stability theorem.

3.1.1. *The singular integral equation.* Let  $a, b \in \mathbb{R}$  with  $a < b$ . On the interval  $\Omega = (a, b)$  we consider the singular integral equation, abbreviated as  $(\lambda\mathbb{I} - A_\Omega)u = f$ ,

$$(3.1) \quad \lambda u(x) - \frac{1}{i\pi} \int_{\Omega} \frac{u(y)}{x-y} dy = f(x), \quad x \in \Omega.$$

The integral is understood in the Cauchy principal value sense. The kernel function  $K(x) = \frac{1}{i\pi x}$  has the Fourier transform

$$\widehat{K}(\xi) = \text{sign } \xi.$$

The operator  $A$  of convolution with  $K$  on  $\mathbb{R}$  is the Hilbert transformation. It satisfies  $A^2 = \mathbb{I}$ , and its spectrum (in a large class of function spaces, for instance  $L^p(\mathbb{R})$  with  $1 < p < \infty$ ) is  $\{-1, 1\}$ , consisting of two eigenvalues of infinite multiplicity.

The finite Hilbert transformation  $A_\Omega$  and its spectral theory are also well studied classical objects, see for example [11]. Here the spectrum depends on the function space; for  $L^p(\Omega)$  it is strongly dependent on  $p$ , but not on  $\Omega$ , as long as  $\Omega$  is a proper subinterval of  $\mathbb{R}$ . For  $p = 2$  one has the following description.

**Lemma 3.1.** *The finite Hilbert transformation  $A_\Omega$  is a bounded selfadjoint operator in  $L^2(\Omega)$ , unitarily equivalent to the operator of multiplication by  $\sigma$  in  $L^2(-1, 1)$  with  $\sigma(\xi) = \xi$ . Both the spectrum  $\text{Sp}(A_\Omega)$  and the closure of the numerical range  $\overline{W}(A_\Omega)$  are equal to  $\mathcal{C} = [-1, 1]$ . For all  $\lambda \in \mathbb{C} \setminus \mathcal{C}$  and any  $f \in L^2(\Omega)$ , the integral equation (3.1) has a unique solution  $u \in L^2(\Omega)$ , and for the resolvent one has in the  $L^2(\Omega)$  operator norm*

$$(3.2) \quad \|(\lambda\mathbb{I} - A_\Omega)^{-1}\| = \text{dist}(\lambda, \mathcal{C})^{-1}.$$

Explicit formulas for the resolvent are known. For the infinite Hilbert transformation this is trivially obtained by algebra:

$$(\lambda\mathbb{I} - A)^{-1} = \frac{1}{\lambda^2 - 1}(\lambda\mathbb{I} + A),$$

and for the finite Hilbert transformation, formulas for the resolvent can be found for example in [16] or [17].

3.1.2. *The discrete system.* We use the notation of Section 1.2 with  $d = 1$ , in particular  $x_m^N = a^N + \frac{m}{N}$  and  $\omega^N = \{m \in \mathbb{Z} \mid x_m^N \in \Omega\}$ . The simple delta-delta discretization of our singular integral equation (3.1) is

$$(3.3) \quad \lambda u_m - \frac{1}{i\pi N} \sum_{n \in \omega^N, m \neq n} \frac{u_n}{x_m^N - x_n^N} = f(x_m^N), \quad (m \in \omega^N).$$

The system matrix  $T^N$  with matrix elements  $\frac{1}{i\pi N} \frac{1}{x_m^N - x_n^N}$  ( $m, n \in \omega^N$ ) is a finite section of the infinite Toeplitz matrix

$$T = \left( \frac{1}{i\pi(m-n)} \right)_{m,n \in \mathbb{Z}} \quad \text{with zero on the diagonal.}$$

The symbol  $F(\tau)$  is now given by the Fourier series

$$(3.4) \quad F(\tau) = \sum_{m \in \mathbb{Z}, m \neq 0} \frac{e^{im\tau}}{i\pi m} = \sum_{m=1}^{\infty} \frac{2 \sin m\tau}{\pi m}, \quad \tau \in Q = [-\pi, \pi].$$

This series converges for all  $t \in Q$  to the well known saw-tooth function

$$(3.5) \quad F(\tau) = \text{sign } \tau - \frac{\tau}{\pi} \quad (\tau \neq 0), \quad F(0) = 0.$$

The range of this function is the interval  $(-1, 1)$ .

Properties of the matrix  $T$  follow immediately from this symbol  $F$  and can be summarized as follows.

**Lemma 3.2.** *The infinite Toeplitz matrix  $T$  defines a bounded selfadjoint operator in  $\ell^2(\mathbb{Z})$ , unitarily equivalent to the operator of multiplication by  $F$  in  $L^2(-\pi, \pi)$  with  $F$  given in (3.5). Both the spectrum  $\text{Sp}(T)$  and the closure of the numerical range  $\overline{W}(T)$  are equal to  $\mathcal{C} = [-1, 1]$ . For all  $\lambda \in \mathbb{C} \setminus \mathcal{C}$  the operator  $\lambda\mathbb{I} - T$  is invertible in  $\ell^2(\mathbb{Z})$ , and for the resolvent one has in the  $\ell^2(\mathbb{Z})$  operator norm*

$$(3.6) \quad \|(\lambda\mathbb{I} - T)^{-1}\| = \text{dist}(\lambda, \mathcal{C})^{-1}.$$

**Corollary 3.3.** *The matrix  $T^N$  of the system (3.3) is selfadjoint with its eigenvalues in  $\mathcal{C} = [-1, 1]$ . For  $\lambda \in \mathbb{C} \setminus \mathcal{C}$ , there is a uniform resolvent estimate in the  $\ell^2$  operator norm*

$$(3.7) \quad \|(\lambda\mathbb{I} - T^N)^{-1}\| \leq \text{dist}(\lambda, \mathcal{C})^{-1}.$$

The converse is also true: If there is a uniform stability estimate

$$\|(\lambda\mathbb{I} - T^N)^{-1}\| \leq C \quad \text{for all } N,$$

then one also has (by a standard Galerkin argument)  $\|(\lambda\mathbb{I} - T)^{-1}\| \leq C$ , hence  $\text{dist}(\lambda, \mathcal{C}) \geq \frac{1}{C}$  and  $\lambda \notin \mathcal{C}$ . Combining this with Lemma 3.1, we obtain the following description of the stability result for our delta-delta discretisation of the finite Hilbert transform.

**Theorem 3.4.** *For  $\lambda \in \mathbb{C}$  the following are equivalent:*

- (i) *The singular integral equation (3.1) has a unique solution  $u \in L^2(\Omega)$  for any  $f \in L^2(\Omega)$ .*
- (ii) *The discretization method (3.3) is stable in the  $\ell^2$  norm.*

(iii)  $\lambda \notin \mathcal{C}$ , where  $\mathcal{C} = [-1, 1]$ .

For such  $\lambda$ , there is an estimate for the operator norms

$$(3.8) \quad \|(\lambda\mathbb{I} - T^N)^{-1}\|_{\mathcal{L}(\ell^2(\omega^N))} \leq \|(\lambda\mathbb{I} - A_\Omega)^{-1}\|_{\mathcal{L}(L^2(\Omega))}.$$

### 3.2. Example 2. Dimension $d = 2$ , kernel $x_1x_2|x|^{-4}$ .

We consider now the simplest higher-dimensional example where in the notation of Section 1.5.2  $d = 2$  and  $p(x) = -\frac{1}{\pi}b_{12}(x)$ , see (1.29). We show that the stability of its delta-delta approximation follows a similar simple pattern as in the previous one-dimensional example, although the proof is non-trivial.

3.2.1. *The singular integral equation.* The kernel and its Fourier transform are

$$(3.9) \quad K(x) = -\frac{x_1x_2}{\pi|x|^4}, \quad \widehat{K}(\xi) = \frac{\xi_1\xi_2}{|\xi|^2}.$$

For  $\Omega \subset \mathbb{R}^2$ , we consider the singular integral equation  $(\lambda\mathbb{I} - A_\Omega)u = f$  as in (1.24)

$$(3.10) \quad \lambda u(x) - \text{p.v.} \int_{\Omega} K(x-y)u(y) dy = f(x).$$

Observing that the range of the function  $\widehat{K}$  is the interval  $[-\frac{1}{2}, \frac{1}{2}]$ , we can formulate the result of Proposition 1.2 as follows

**Lemma 3.5.** *Let  $\mathcal{C} = [-\frac{1}{2}, \frac{1}{2}]$ . For  $\Omega = \mathbb{R}^2$ , both the spectrum  $\text{Sp}(A_\Omega)$  and the closure of the numerical range  $\overline{W}(A_\Omega)$  in  $L^2(\Omega)$  are equal to  $\mathcal{C}$ . For any open subset  $\Omega \subset \mathbb{R}^2$ , the closure of the numerical range in  $L^2(\Omega)$  satisfies  $\overline{W}(A_\Omega) \subset \mathcal{C}$ , and there is a resolvent estimate in the  $L^2(\Omega)$  operator norm*

$$(3.11) \quad \|(\lambda\mathbb{I} - A_\Omega)^{-1}\| \leq \text{dist}(\lambda, \mathcal{C})^{-1}.$$

3.2.2. *The discrete system.* Let now  $\Omega$  be a bounded domain in  $\mathbb{R}^2$ . In the notation of Section 1.6 with  $d = 2$ , the regular grid consists of the points  $x_m^N = a^N + \frac{m}{N}$ , indexed by  $\omega^N = \{m \in \mathbb{Z}^2 \mid x_m^N \in \Omega\}$ . The simple delta-delta discretization of our singular integral equation (3.10) is

$$(3.12) \quad \lambda u_m + \frac{1}{\pi}N^{-2} \sum_{n \in \omega^N, n \neq m} \frac{(x_{m,1}^N - x_{n,1}^N)(x_{m,2}^N - x_{n,2}^N)}{|x_m^N - x_n^N|^4} u_n = f(x_m^N), \quad (m \in \omega^N).$$

The system matrix  $T^N$  is now a finite section of the infinite Toeplitz matrix

$$T = -\frac{1}{\pi} \left( \frac{(m_1 - n_1)(m_2 - n_2)}{|m - n|^4} \right)_{m, n \in \mathbb{Z}^2} \quad \text{with zero on the diagonal.}$$

Its symbol is therefore given by the Fourier series for  $\tau \in Q = [-\pi, \pi]^2$

$$(3.13) \quad F(\tau) = - \sum_{m \in \mathbb{Z}^2, m \neq 0} \frac{m_1 m_2}{\pi|m|^4} e^{im \cdot \tau} = \frac{4}{\pi} \sum_{m_1, m_2=1}^{\infty} \frac{m_1 m_2}{(m_1^2 + m_2^2)^2} \sin(m_1 \tau_1) \sin(m_2 \tau_2).$$

Whereas we do not know an explicit closed form expression for this function, we know from the results of Section 2.2 using Ewald's method that it is bounded and that it can be written as in equation (2.13)

$$(3.14) \quad F(\tau) = \widehat{K}(\tau) + F_0(\tau) \quad \text{where } F_0 \text{ is analytic in } Q \text{ and } F_0(0) = 0.$$

In addition, we know from (3.13) that  $F$  vanishes on the boundary of  $Q$ , hence

$$(3.15) \quad F_0(\tau) = -\widehat{K}(\tau) \quad (\tau \in \partial Q).$$

In the previous example, we used the explicit expression of  $F(\tau)$  for finding the range of  $F$ . In fact, the function  $F_0$  in that case was just the linear interpolation between the two values of the symbol  $\widehat{K}$  on  $\partial Q$ , which implied that the closed convex hull of  $\text{im}(F)$  was the same as the convex hull of  $\text{im}(\widehat{K})$ . In the present case, we do not have a simple formula, but we can still prove that the conclusion is true.

**Lemma 3.6.** *Let  $F(\tau)$  be as defined in (3.13). Then for any  $\tau \in Q$*

$$(3.16) \quad F(\tau) \in \mathcal{C} = \left[-\frac{1}{2}, \frac{1}{2}\right].$$

The proof is not obvious, although the claim is numerically evident if we compute  $F$  using Ewald's method and plot its graph, see the contour plot in Figure 1.

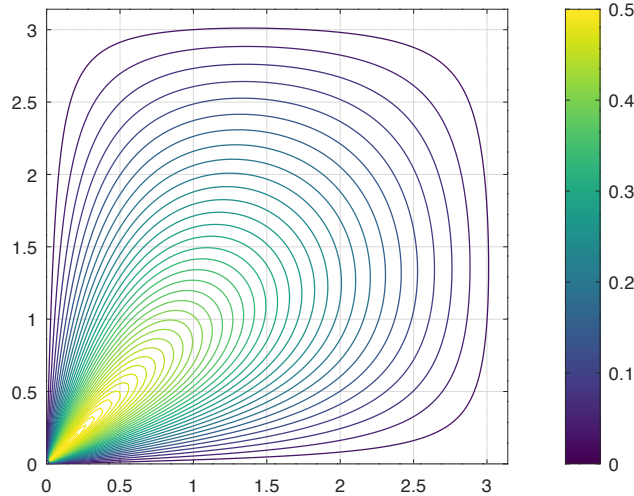


FIGURE 1. Contour plot of  $F(\xi)$  on the quarter square  $Q_{++}$ , Example 2.

Before we give the proof, let us draw the conclusion for the stability of the numerical scheme (3.12).

**Corollary 3.7.** *Let  $\mathcal{C} = [-\frac{1}{2}, \frac{1}{2}]$  and  $\lambda \in \mathbb{C} \setminus \mathcal{C}$ . Then for any  $N$  the linear system (3.12) has a unique solution, and there is a uniform resolvent estimate*

$$(3.17) \quad \|(\lambda \mathbb{I} - T^N)^{-1}\|_{\mathcal{L}(\ell^2(\Omega^N))} \leq \text{dist}(\lambda, \mathcal{C})^{-1}.$$

*Proof.* For symmetry reasons, it is sufficient to prove (3.16) for  $\tau \in Q_{++} = (0, \pi)^2$ . For the proof of Lemma 3.6, we will show the following:

$$(3.18) \quad \text{For any } \xi \in Q_{++}, \quad F(\xi) \geq 0 \quad \text{and} \quad F_0(\xi) \leq 0.$$

This implies  $0 \leq F(\xi) \leq \widehat{K}(\xi) \leq \frac{1}{2}$ , hence (3.16).

We use the integral representations from Proposition 2.5

$$(3.19) \quad F_0(\xi) = \int_0^\infty H_0(\xi, s) ds, \quad \widehat{K}(\xi) = \int_0^\infty H_1(\xi, s) ds, \quad F(\xi) = \int_0^\infty H(\xi, s) ds.$$

According to (2.19),  $H_0(\xi, s) = H(\xi, s) - H_1(\xi, s)$  with

$$(3.20) \quad H_1(\xi, s) = \frac{\xi_1 \xi_2}{4s^2} e^{-\frac{|\xi|^2}{4s}}, \quad H(\xi, s) = \sum_{n \in \mathbb{Z}^2} H_1(\xi + 2\pi n, s).$$

Let  $0 \leq \epsilon < T$  and  $\Sigma_\epsilon^T = Q_{++} \times (\epsilon, T)$ . In  $\Sigma_\epsilon^T$ , we want to use the maximum principle for the heat equation (see Lemma 2.6) for the functions  $\tilde{H}_0(\xi, s) = s^{-1}H_0(\xi, s)$  and  $\tilde{H}(\xi, s) = s^{-1}H(\xi, s)$ .

Since  $\tilde{H}_1(\xi, s) = s^{-1}H_1(\xi, s)$  is continuous for  $(\xi, s) \in \mathbb{R}^2 \times [0, \infty) \setminus \{0, 0\}$  and the Poisson series

$$\tilde{H}_0(\xi, s) = \sum_{n \in \mathbb{Z}^2, n \neq 0} \tilde{H}_1(\xi + 2\pi n, s)$$

converges uniformly for  $(\xi, s) \in \overline{\Sigma_0^T}$  for all  $T > 0$ , we see that  $\tilde{H}_0$  is continuous in  $\overline{\Sigma_0^T}$  with initial value  $\tilde{H}_0(\xi, 0) = 0$ . On the lateral boundary we use the Fourier representation (see (2.16))

$$\tilde{H}(\xi, s) = \frac{4}{\pi} \sum_{m_1, m_2=1}^{\infty} m_1 m_2 e^{-|m|^2 s} \sin(m_1 \xi_1) \sin(m_2 \xi_2).$$

If  $\xi_1$  or  $\xi_2$  is in  $\{0, \pi\}$ , this implies that  $\tilde{H} = 0$  and therefore

$$\tilde{H}_0(\xi, s) = -\tilde{H}_1(\xi, s) \leq 0 \quad \text{for } (\xi, s) \in \partial Q_{++} \times (0, T].$$

According to Lemma 2.6,  $\tilde{H}_0$  satisfies the heat equation  $(\partial_s - \Delta_\xi)\tilde{H}_0 = 0$  in  $\Sigma_0^T$ . Thus we can apply the maximum principle to  $\tilde{H}_0$  and obtain  $\tilde{H}_0(\xi, s) \leq 0$  in  $\Sigma_0^T$ , hence also  $H_0(\xi, s) \leq 0$ . Integrating over  $s \in (0, \infty)$  yields

$$F_0(\xi) \leq 0 \quad \text{for } \xi \in Q_{++}.$$

For  $\tilde{H}$ , we cannot apply the maximum principle directly in  $\Sigma_0^T$ , because  $\tilde{H}$  is not continuous at  $(0, 0) \in \overline{\Sigma_0^T}$ , but we can apply it in  $\Sigma_\epsilon^T$  for any  $0 < \epsilon < T$ . On the lateral boundary,  $\tilde{H}$  vanishes as seen above, and for the initial value at  $s = \epsilon$  we have

$$\tilde{H}(\xi, \epsilon) = \tilde{H}_0(\xi, \epsilon) + \tilde{H}_1(\xi, \epsilon) \geq \tilde{H}_0(\xi, \epsilon) \geq \delta(\epsilon)$$

with  $\delta(\epsilon) = \inf_{\xi \in Q_{++}} \tilde{H}_0(\xi, \epsilon)$ . Hence by the maximum principle, in  $\overline{\Sigma_\epsilon^T}$  we have

$$\tilde{H}(\xi, s) \geq \min\{0, \delta(\epsilon)\}.$$

Now, as we have seen above,  $\tilde{H}_0(\cdot, s)$  tends to 0 uniformly as  $s \rightarrow 0$ , hence  $\delta(\epsilon) \rightarrow 0$  as  $\epsilon \rightarrow 0$ , which implies  $\tilde{H}(\xi, s) \geq 0$  for any  $s > 0$  and  $\xi \in Q_{++}$ . After integrating over  $s$ , we finally get  $F(\xi) \geq 0$  for  $\xi \in Q_{++}$ , and the proof of the Lemma is complete.  $\square$

*Remark 3.8.* In conclusion, for this example we find the same “ideal” stability estimate as in the previous one-dimensional example.

### 3.3. Example 3. Dimension $d = 2$ , kernel $(x_1^2 - x_2^2)|x|^{-4}$ .

We consider another two-dimensional example where in the notation of Section 1.5.2  $d = 2$  and  $p(x) = -\frac{1}{2\pi}a_{12}(x)$ , see (1.29). We show that the complement of the stability zone in this case is strictly larger than the image of the symbol of the integral operator.

3.3.1. *The singular integral equation.* We use the same notation for analogous objects as in the preceding example. Therefore in this section, the letters  $K, \widehat{K}, T$  etc. are redefined to have new meanings. The kernel and its Fourier transform are now

$$(3.21) \quad K(x) = \frac{x_2^2 - x_1^2}{2\pi|x|^4}, \quad \widehat{K}(\xi) = \frac{\xi_1^2 - \xi_2^2}{2|\xi|^2} = \frac{\xi_1^2}{|\xi|^2} - \frac{1}{2}.$$

The normalization is chosen so that the range of the function  $\widehat{K}$  is again the interval  $[-\frac{1}{2}, \frac{1}{2}]$ .

In fact, this kernel is the same as in the previous example (3.9) after a  $45^\circ$  rotation of the coordinate system. Therefore if we write the singular integral equation as in (3.10), we can copy verbatim the statement of the previous example concerning the numerical range of the integral operator  $A_\Omega$  (see Lemma 3.5) and the resolvent estimate (3.11).

**Lemma 3.9.** *Lemma 3.5 is true for the singular integral equation (3.10) defined with the kernel (3.21).*

3.3.2. *The discrete system.* To the delta-delta discretization

$$(3.22) \quad \lambda u_m - N^{-2} \sum_{n \in \omega^N, m \neq n} K(x_m^N - x_n^N) u_n = f(x_m^N) \quad (m \in \omega^N)$$

corresponds the finite section  $T^N$  of the infinite Toeplitz matrix

$$T = \frac{1}{2\pi} \left( \frac{(m_2 - n_2)^2 - (m_1 - n_1)^2}{|m - n|^4} \right)_{m, n \in \mathbb{Z}^2} \quad \text{with zero on the diagonal.}$$

The numerical symbol (symbol of  $T$ ) is now defined as

$$(3.23) \quad F(\tau) = \sum_{m \in \mathbb{Z}^2, m \neq 0} \frac{m_2^2 - m_1^2}{2\pi|m|^4} e^{im \cdot \tau}.$$

**Lemma 3.10.** *Let*

$$\Lambda_0 = \frac{\Gamma(\frac{1}{4})^4}{32\pi^2} = 0.5471\dots$$

*Let  $F(\tau)$  be as defined in (3.23). Then there exists  $\Lambda_+ \geq \Lambda_0$  such that  $F(Q) = \mathcal{C} = [-\Lambda_+, \Lambda_+]$ .*

**Conjecture 3.11.** *Numerical evidence suggests equality*

$$(3.24) \quad \Lambda_+ = \Lambda_0.$$

*Proof of Lemma 3.10.* The function  $F$  is odd with respect to permutation of  $\xi_1$  and  $\xi_2$ . The decomposition  $F = F_0 + \widehat{K}$  with  $F_0$  continuous on  $Q$  implies that  $F$  takes its maximum  $\Lambda_+$  on  $Q$ . Therefore its image  $F(Q)$  is a closed symmetric interval  $\mathcal{C} = [-\Lambda_+, \Lambda_+]$ . We are going to show that

$$(3.25) \quad F(\pi, 0) = \Lambda_0.$$

The conjecture (3.24) then corresponds to the claim that  $F$  attains its maximum on  $Q$  in the point  $\tau = (\pi, 0)$ .

To prove (3.25), we first transform the slowly converging double Fourier series

$$(3.26) \quad F(\pi, 0) = \sum_{m \in \mathbb{Z}^2, m \neq 0} (-1)^{m_1} \frac{m_2^2 - m_1^2}{2\pi(m_1^2 + m_2^2)^2}$$

into a rapidly convergent single series. One way to get this is to start with the Poisson summation formula applied to the function  $f(x) = (x - iy)^{-1}$  whose Fourier transform is  $\widehat{f}(\xi) = 2\pi i \mathbb{1}_+(\xi) e^{-y\xi}$  for  $y > 0$ . The result is then valid for all  $y \neq 0$ . It can be written for  $t \in [-\pi, \pi]$  as

$$(3.27) \quad \sum_{n \in \mathbb{Z}} \frac{e^{int}}{n - iy} = i\pi \frac{e^{y\sigma(t)}}{\sinh(\pi y)} \quad \text{with } \sigma(t) = -t + \pi \operatorname{sign} t.$$

Remark: Euler's formula (3.5) is a simple consequence of this.

Taking the derivative with respect to  $y$  and subtracting the formulas for  $y$  and  $-y$  leads to

$$(3.28) \quad \sum_{n \in \mathbb{Z}} \frac{n^2 - y^2}{(n^2 + y^2)^2} e^{int} = \pi \frac{\sigma(t) \sinh(\sigma(t)y) \sinh(\pi y) - \pi \cosh(\sigma(t)y) \cosh(\pi y)}{\sinh^2 \pi y}.$$

This can be used to reduce the double Fourier series for  $F(\xi)$  to a single rapidly convergent Fourier series. We are here only interested in the limit  $t \rightarrow 0$ :

$$(3.29) \quad \sum_{n \in \mathbb{Z}} \frac{n^2 - y^2}{(n^2 + y^2)^2} = \frac{-\pi^2}{\sinh^2 \pi y}.$$

Hence, by decomposing the double sum  $\sum_{m \in \mathbb{Z}^2 \setminus \{0\}}$  as  $\sum_{m_1=0, m_2 \in \mathbb{Z} \setminus \{0\}} + \sum_{m_1 \in \mathbb{Z} \setminus \{0\}, m_2 \in \mathbb{Z}}$  and using  $\sum_{n=1}^{\infty} \frac{1}{n^2} = \frac{\pi^2}{6}$  (which can also be obtained from (3.29) by looking at the pole in  $y = 0$ ) we finally get

$$(3.30) \quad \Lambda_0 = \sum_{m \in \mathbb{Z}^2, m \neq 0} (-1)^{m_1} \frac{m_2^2 - m_1^2}{2\pi(m_1^2 + m_2^2)^2} = \frac{\pi}{6} - \sum_{n=1}^{\infty} \frac{(-1)^n \pi}{\sinh^2 \pi n}.$$

This series converges rapidly, with 5 terms giving 15 significant digits:  $\Lambda_0 = 0.547109903806619\dots$ . The series is covered by the formulas for  $\operatorname{IX}_s$  with  $s = 2$  in [21]. The explicit expression for  $\Lambda_0$  given in the Lemma can be deduced from this.  $\square$

*Remark 3.12.* The conformal radius (or logarithmic capacity) of the unit square is known to be [12, Tables]

$$R_{\square} = \frac{\Gamma(\frac{1}{4})^2}{4\pi^{\frac{3}{2}}}.$$

This implies the remarkable relation

$$(3.31) \quad \pi R_{\square}^2 = 2\Lambda_0.$$



The conjecture that  $\Lambda_+ = \Lambda_0$  is clearly supported by numerical evidence. Here are the results of two different approaches for the approximation of  $\Lambda_+$ :

In Table 1 we approximate the numerical symbol  $F$  from (3.23) using the Ewald method (2.12) from Proposition 2.2.

$$(3.32) \quad F(\tau) \approx \sum_{|m_1|, |m_2| \leq M, m \neq 0} K(m) \Gamma(2, \pi |m|^2) e^{im \cdot \tau} + \sum_{|n_1|, |n_2| \leq M} \widehat{K}(\tau + 2\pi n) e^{-\frac{|\tau + 2\pi n|^2}{4\pi}}.$$

We take the maximum of  $F(\tau)$  over a regular  $N \times N$  grid discretizing the period square  $Q = [-\pi, \pi]^2$ . Results are shown for  $N = 1001$ , so that the point  $(\pi, 0)$  is included. One sees the rapid convergence of the sums in the Ewald method.

$M$	Maximum	diff with $\Lambda_0$
1	0.5466820485568409	-0.00043
2	0.5471099022284376	-1.578e-9
3	0.5471099038066192	1.11e-16
4	0.5471099038066192	1.11e-16

TABLE 1. Computation of  $\Lambda_+$

In Table 2 we show the maximum eigenvalue of the matrix  $T^N$  where  $\Omega$  is the unit square, together with an extrapolated value and its difference with  $\Lambda_0$ .

$N$	$\lambda_{\max}(T^N)$	extrap.	diff with $\Lambda_0$
16	0.541802946417726		
24	0.544571778645890		
36	0.545922219922679	0.547207966733364	9.81e-5
54	0.546562896841136	0.547141211191569	3.13e-5
81	0.546860792009930	0.547119678405314	9.77e-6

TABLE 2. Computation of  $\max(\text{Sp}(T^N))$ , Example 3

For comparison, we show in Table 3 the analogous computations for the matrices from Example 2, where  $\Lambda_+ = 0.5$ .

In the previous Example 2, we were able to prove the equation  $\Lambda_+ = 0.5$  using an argument involving the maximum principle for the heat equation, see the proof of Lemma 3.6, in particular (3.18). While we have no proof for the equation  $\Lambda_+ = \Lambda_0$  here, it is possible to use an analogous argument to obtain an upper bound for  $\Lambda_+$ . The square  $Q_{++}$  (of area  $\pi^2$ ) of the previous example now has to be turned by  $45^\circ$  and to be replaced by the lozenge (a square of area  $2\pi^2$ )

$$Q_\diamond = \{\xi \in \mathbb{R}^2 \mid 0 < \xi_1 < 2\pi; |\xi_2| < \min\{\xi_1, 2\pi - \xi_1\}\}.$$

$N$	$\lambda_{\max}(T^N)$	extrap.	diff with 0.5
16	0.4299869696672885		
24	0.4526591158216325		
36	0.4683227545642122	0.5033301483277116	3.33e-3
54	0.4789372344435390	0.5012512843991882	1.25e-3
81	0.4860451011088278	0.5004526691660266	4.53e-4

TABLE 3. Computation of  $\max(\text{Sp}(T^N))$ , Example 2

Then one can see again that  $F(\xi) = 0$  on  $\partial Q_\diamond$ . But now  $\partial Q_\diamond$  contains two points of discontinuity of  $F$ , the origin  $(0, 0)$  and the point  $(2\pi, 0)$ . Therefore the decomposition  $F = F_0 + \widehat{K}$  has to be refined into

$$F(\xi) = F_{00}(\xi) + \widehat{K}(\xi) + \widehat{K}(\xi - (2\pi, 0)).$$

The function  $F_{00}$  defined by this will then be continuous on the closure of  $Q_\diamond$ . Now one can use the integral representation from Proposition 2.5 similarly to (3.19) and use the maximum principle for the heat equation as before to conclude that

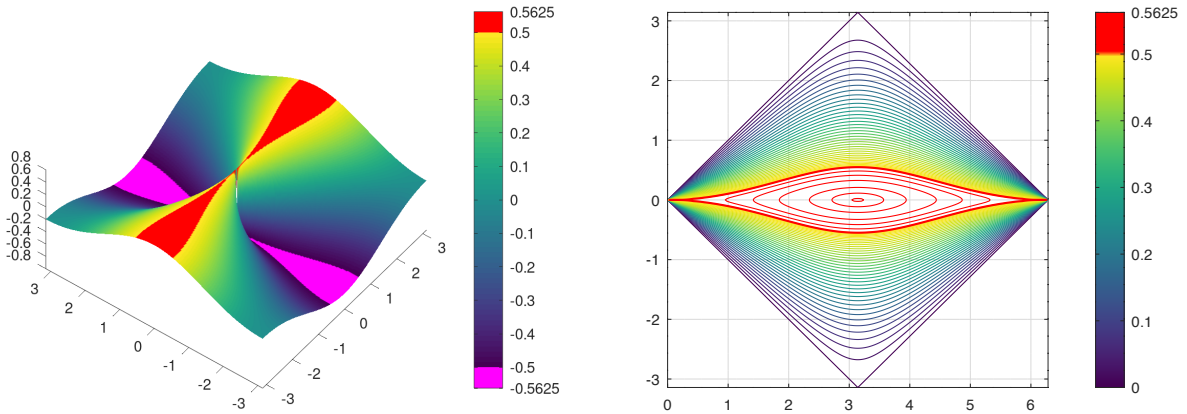
$$\text{For any } \xi \in Q_\diamond, \quad F(\xi) \geq 0 \text{ and } F_{00}(\xi) \leq 0.$$

This implies  $0 \leq F(\xi) \leq \widehat{K}(\xi) + \widehat{K}(\xi - (2\pi, 0))$  in  $Q_\diamond$ , and hence by taking the maximum,

$$(3.33) \quad \Lambda_+ \leq 1.$$

Unfortunately, this estimate is much less sharp than the estimate by  $\frac{1}{2}$  in the previous example.

To illustrate the behavior of the numerical symbol  $F(\xi)$ , we present in Figure 2 a surface graph of  $F$  on the square  $Q$  and a contour plot on the lozenge  $Q_\diamond$ . The parts exceeding the range of the symbol  $\widehat{K}$  are indicated in bright red hues. The maximum at the midpoint  $(\pi, 0)$  of  $Q_\diamond$  is in clear evidence.

FIGURE 2. Numerical symbol for Example 3. Left:  $F$  on  $Q$ , right:  $F$  on  $Q_\diamond$ .

Let us summarize the stability result obtained for this example.

**Corollary 3.13.** *Let  $\mathcal{C} = [-\Lambda_+, \Lambda_+]$  and  $\lambda \in \mathbb{C} \setminus \mathcal{C}$ . Then for any  $N$  the linear system (3.12) has a unique solution, and there is a uniform resolvent estimate*

$$(3.34) \quad \|(\lambda \mathbb{I} - T^N)^{-1}\|_{\mathcal{L}(\ell^2(\Omega^N))} \leq \text{dist}(\lambda, \mathcal{C})^{-1}.$$

*For  $\lambda \in [-\Lambda_+, -\frac{1}{2}] \cup (\frac{1}{2}, \Lambda_+]$  the integral equation  $(\lambda \mathbb{I} - A_\Omega)u = f$  with kernel (3.21) is well-posed in  $L^2(\Omega)$ , but the corresponding delta-delta approximation scheme (3.22) is unstable.*

#### 3.4. Example 4. Dimension $d = 2$ , kernel $(x_1 + ix_2)^2|x|^{-4}$ .

Let  $d = 2$  and  $p(x) = -\frac{1}{\pi}(a_{12}(x) + 2ib_{12}(x))$ . The corresponding kernel and its Fourier transform are

$$(3.35) \quad K(x) = \frac{x_2^2 - x_1^2 - 2ix_1x_2}{\pi|x|^4}, \quad \widehat{K}(\xi) = \frac{(\xi_1 + i\xi_2)^2}{|\xi_1 + i\xi_2|^2}.$$

The normalization is chosen so that  $|\widehat{K}(\xi)| = 1$  for  $\xi \in \mathbb{R}^2$ . We include this example, which has features combining those of the two preceding examples, mainly for purposes of illustration. Because the singular integral operator and the system matrices of the corresponding delta-delta discretization in this case are non-selfadjoint, we expect to see less trivial relations between spectra and numerical ranges than in the selfadjoint case.

It is obvious from the definition (3.35) that the spectrum  $\text{Sp}(A)$  of the operator of convolution with  $K$  in  $L^2(\mathbb{R}^2)$  is the unit circle  $\{\xi \in \mathbb{C} \mid |\xi| = 1\}$  and that its numerical range is the unit disk. Whereas we do not know the spectrum  $\text{Sp}(A_\Omega)$  for a bounded domain  $\Omega \subset \mathbb{R}^2$ , the numerical range is still the unit disk, compare (1.26),

$$(3.36) \quad \text{Sp}(A_\Omega) \subset W(A_\Omega) = W(A) = \{\xi \in \mathbb{C} \mid |\xi| \leq 1\}.$$

For the system matrices  $T^N$  of the delta-delta discretization scheme, Theorem 1.6 and Lemma 2.1 provide the following relations.

$$(3.37) \quad \text{Sp}(T^N) \subset W(T^N) \subset \overline{W(T)} = \overline{\text{conv}} \bigcup_{M \in \mathbb{N}} W(T^M) \quad \text{and} \quad W(A) \subset \overline{W(T)}.$$

In Figure 3 we show for the case of a square domain  $\Omega$  and two values of  $N$  the spectrum  $\text{Sp}(T^N)$  (red points), the boundary of the numerical range  $W(T^N)$  (red line), and the unit circle, which is the boundary of  $W(A)$  (green line). We can see the inclusions from (3.37) between  $\text{Sp}(T^N)$  and  $W(T^N)$ , and we can perceive the asymptotic inclusion of  $W(A)$  in  $W(T^N)$  as  $N$  tends to infinity.

We can also see that the eigenvalues of the matrices, in contrast to the numerical range, will not fill the whole unit disk asymptotically. On the other hand, we clearly see the overshoot  $W(T) \setminus W(A_\Omega)$ , that is the region of  $\lambda$  where the volume integral equation is uniquely solvable and the operator  $\lambda \mathbb{I} - A_\Omega$  is sectorial, so that every  $L^2$ -conforming Galerkin method would converge, whereas the delta-delta scheme is unstable. It appears that the limits for the real part of this overshoot are the same (scaled by a factor 2) as in the previous example, that is  $\pm 2\Lambda_0$  with  $\Lambda_0$  defined in Lemma 3.10.

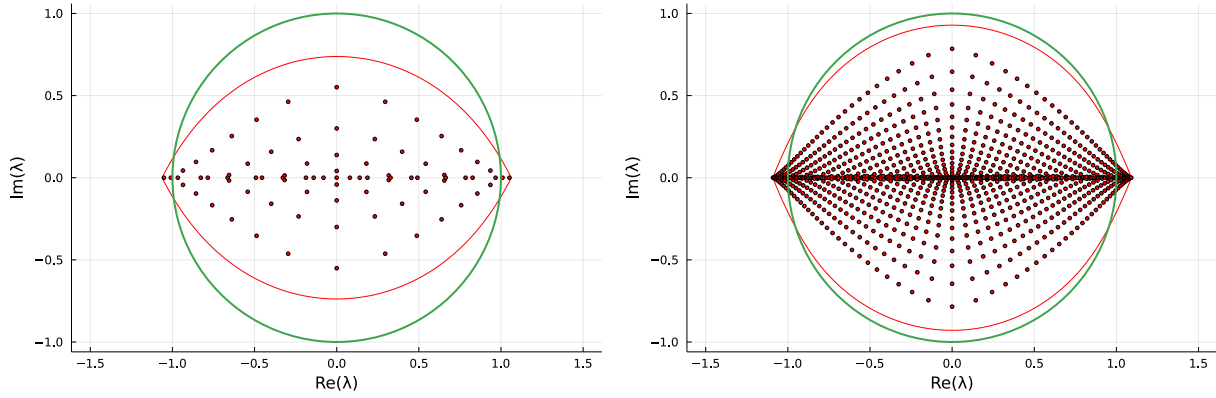


FIGURE 3. Spectrum and numerical range. Left:  $N = 8$ , right:  $N = 32$ .

### 3.5. Example 5. Dimension $d \geq 2$ . Volume Integral Equation for the Quasi-static Maxwell system.

In the quasi-static Maxwell volume integral equation (see Section 1.2), the right hand side and the solution are  $\mathbb{C}^d$ -valued functions, and the singular integral operator is defined as the matrix of second distributional derivatives of the convolution with the free-space Green function  $g$  for the Laplace operator, see equation (1.5). If we call this operator  $A^0$ , then it is not the same as the operator  $A$  defined by the Cauchy principal value of the integral with the same kernel, but there is a simple relation: Let

$$(3.38) \quad K(x) = -D^2g(x) = \left( -\partial_i\partial_jg(x) \right)_{i,j \in \{1, \dots, d\}} \quad \text{for } x \neq 0.$$

Then

$$(3.39) \quad A^0u(x) = -\nabla \operatorname{div} \int_{\mathbb{R}^d} g(x-y)u(y)dy = \text{p.v.} \int_{\mathbb{R}^d} K(x-y)u(y)dy + \frac{1}{d}u(x) = (A + \frac{1}{d}\mathbb{I})u(x).$$

This is most easily seen by first using the symmetries of the kernel with respect to reflections at coordinate axes and permutations of the variables in order to deduce that the distribution kernel of  $A^0 - A$  must be a scalar multiple of the  $d \times d$  identity matrix  $\mathbb{I}_d$  times the Dirac distribution  $\delta_0$ , and then determining this multiple by taking traces:  $\operatorname{tr}(-D^2g) = -\Delta g = \delta_0 = \operatorname{tr}(\frac{1}{d}\mathbb{I}_d\delta_0)$ .

3.5.1. *The singular integral equation.* We consider the strongly singular integral equation, still written as  $(\lambda\mathbb{I} - A_\Omega)u = f$ ,

$$(3.40) \quad \lambda u(x) - \text{p.v.} \int_{\Omega} K(x-y)u(y)dy = f(x) \quad \text{with } K \text{ given in (3.38)}.$$

The function space is now  $L^2(\Omega; \mathbb{C}^d)$ .

Let us note the explicit form of the kernel, valid in any dimension  $d \geq 2$ , where we consider points in  $\mathbb{R}^d$  as column vectors,

$$(3.41) \quad K(x) = -\frac{1}{\nu_d} \left( x x^\top - \frac{1}{d}\mathbb{I}_d|x|^2 \right) |x|^{-d-2}, \quad \text{with } \nu_d = \frac{2\pi^{\frac{d}{2}}}{d\Gamma(\frac{d}{2})}.$$

The simplest way to see this is to first look at the symbol of the operator. For this we employ  $d$ -dimensional Fourier transformation and use the fact that  $\widehat{g}(\xi) = |\xi|^{-2}$ , hence

$$(3.42) \quad \mathcal{F}(-D^2g)(\xi) = \frac{\xi \xi^\top}{|\xi|^2} \quad \text{and} \quad \widehat{K}(\xi) = \frac{\xi \xi^\top - \frac{1}{d}\mathbb{I}_d|\xi|^2}{|\xi|^2}.$$

We check that  $\text{tr } \widehat{K} = 0$  and that  $\widehat{K}$  satisfies the spherical cancellation condition. Indeed, in the notation of Lemma 1.5, the (matrix-valued) polynomial  $p(\xi)$  is given by

$$(3.43) \quad p(\xi) = -\frac{1}{\nu_d}(\xi \xi^\top - \frac{1}{d}\mathbb{I}_d|\xi|^2).$$

Thus the off-diagonal elements of the matrix  $p(x)$  are given by

$$-\frac{1}{\nu_d}x_jx_k = -\frac{b_{jk}(x)}{\nu_d},$$

and the diagonal elements by

$$-\frac{1}{\nu_d}(x_k^2 - \frac{1}{d}|x|^2) = \frac{1}{d\nu_d} \sum_{j=1}^d a_{jk}(x),$$

compare Remark 2.4.

From our formulas of Section 1.5.1 we find the explicit form (3.41) for our kernel. For  $d = 2$ , we have  $\nu_d = \pi$ , and we recognize the kernels studied in the Examples 2 and 3.

The matrix  $\frac{\xi \xi^\top}{|\xi|^2}$  is an orthogonal projection matrix, hence its numerical range is the interval  $[0, 1]$ . Therefore  $W(\widehat{K}(\xi)) = [-\frac{1}{d}, 1 - \frac{1}{d}]$  for any  $\xi \neq 0$ . We immediately get the following instance of Proposition 1.2.

**Lemma 3.14.** *Let  $\mathcal{C} = [-\frac{1}{d}, 1 - \frac{1}{d}]$ . Then for all  $\lambda \notin \mathcal{C}$  and any  $f \in L^2(\Omega; \mathbb{C}^d)$ , the integral equation (3.40) has a unique solution  $u \in L^2(\Omega; \mathbb{C}^d)$ , and there is a resolvent estimate in the  $L^2(\Omega; \mathbb{C}^d)$  operator norm*

$$(3.44) \quad \|(\lambda\mathbb{I} - A_\Omega)^{-1}\| \leq \text{dist}(\lambda, \mathcal{C})^{-1}.$$

**3.5.2. The discrete system.** With the  $d \times d$  matrix-valued kernel  $K$  and vector-valued functions  $u$  and  $f$ , we can write the delta-delta discretization  $(\lambda\mathbb{I} - T^N)U = F$  of the integral equation (3.40) in the same form as in the scalar case

$$(3.45) \quad \lambda u_m - N^{-d} \sum_{n \in \omega^N, m \neq n} K(x_m^N - x_n^N)u_n = f(x_m^N), \quad (m \in \omega^N),$$

where now the system matrix  $T^N$  is of size  $d|\omega^N| \times d|\omega^N|$  and is considered as a linear operator in  $\ell^2(\omega^N; \mathbb{C}^d)$ .

We recall the discussion of matrix-valued kernels in Section 2.4 above, in particular the properties of the numerical range stated in Lemma 2.7.

The basic stability estimate follows.

**Proposition 3.15.** *Let  $K$  be the kernel defined in (3.38), (3.41). Then there exist  $\Lambda_-^{(d)}, \Lambda_+^{(d)} \in \mathbb{R}$  with*

$$(3.46) \quad \Lambda_-^{(d)} \leq -\frac{1}{d}, \quad \Lambda_+^{(d)} \geq 1 - \frac{1}{d}$$

such that the following holds.

(i) For  $\tau \in Q = [-\pi, \pi]^d$ ,  $\tau \neq 0$ ,  $F(\tau)$  is a real symmetric matrix with eigenvalues contained in the interval  ${}^c\mathcal{C} = [\Lambda_-^{(d)}, \Lambda_+^{(d)}]$ ,

$$(3.47) \quad \Lambda_-^{(d)} = \inf_{\tau \in Q} \min(\text{Sp}(F(\tau))), \quad \Lambda_+^{(d)} = \sup_{\tau \in Q} \max(\text{Sp}(F(\tau))).$$

(ii) For any  $N$ , the numerical range  $W(T^N)$  is contained in  $W(T) = {}^c\mathcal{C}$ .

$$(3.48) \quad \Lambda_-^{(d)} = \inf_{N \in \mathbb{N}} \min(\text{Sp}(T^N)), \quad \Lambda_+^{(d)} = \sup_{N \in \mathbb{N}} \max(\text{Sp}(T^N)).$$

(iii) The delta-delta scheme (3.45) is stable if and only if  $\lambda \in \mathbb{C} \setminus {}^c\mathcal{C}$ , and one has the stability estimate in the  $\ell^2(\omega^N; \mathbb{C}^d)$  operator norm

$$(3.49) \quad \|(\lambda \mathbb{I} - T^N)^{-1}\| \leq \text{dist}(\lambda, {}^c\mathcal{C})^{-1}.$$

*Proof.* The matrix  $K(x)$  is symmetric for  $x \neq 0$ , implying that also  $F(\tau)$  is a symmetric matrix for  $\tau \neq 0$ . The symmetry  $K(-x) = K(x)$  implies that the matrix elements of  $F(\tau)$  are real. Therefore the numerical range of  $F(\tau)$  is the interval  $[\lambda_-(\tau), \lambda_+(\tau)]$ , where

$$\lambda_-(\tau) = \min(\text{Sp}(F(\tau))), \quad \lambda_+(\tau) = \max(\text{Sp}(F(\tau))).$$

This justifies (3.47). All the other statements of the proposition are instances of the statements of Section 1.6, in particular Theorem 1.6, and their proofs in Section 2, based on Ewald's method.  $\square$

What remains is to get information on the numbers  $\Lambda_{\pm}^{(d)}$  and to see whether the inequalities (3.46) are strict. In that case, for  $\lambda \in [\Lambda_-^{(d)}, -\frac{1}{d}] \cup (1 - \frac{1}{d}, \Lambda_+^{(d)}]$ , the integral equation is well-posed in  $L^2(\Omega; \mathbb{C}^d)$ , but the delta-delta discretization scheme is unstable in  $\ell^2(\Omega^N; \mathbb{C}^d)$ .

We will discuss this for the practically relevant cases  $d = 2$ , where we get rather precise information, and  $d = 3$ , which is the most important case because of its relevance for the DDA method in computational electromagnetics.

3.5.3. *Dimension  $d = 2$ .* Here the numerical symbol has the form

$$F(\tau) = \begin{pmatrix} a(\tau) & b(\tau) \\ b(\tau) & -a(\tau) \end{pmatrix}$$

with real-valued functions  $a$  and  $b$ . The eigenvalues are  $\lambda_{\pm}(\tau) = \pm \sqrt{a(\tau)^2 + b(\tau)^2}$ , implying  $\Lambda_-^{(d)} = -\Lambda_+^{(d)}$ . The functions  $a$  and  $b$  have been studied in the previous examples,  $a$  in Example 3 and  $b$  in Example 2.

In particular,  $b(\tau) = 0$  for  $\tau \in \partial Q$ , and therefore

$$(3.50) \quad \text{for } \tau = (\pi, 0), \quad \lambda_+(\tau) = a(\tau) = \Lambda_0$$

with the number  $\Lambda_0 = 0.5471\dots$  encountered in Example 3, Lemma 3.10. This implies

$$\Lambda_+^{(2)} \geq \Lambda_0,$$

and we are in the same situation as in Example 3: Strong numerical evidence suggests that the function  $\lambda_+$  attains its maximum in the point  $\tau = (\pi, 0)$  and therefore  $\Lambda_+^{(2)} = \Lambda_0$ , but we do not have a formal proof for this. The positive eigenvalue  $\lambda_+$  is plotted in Figure 4.

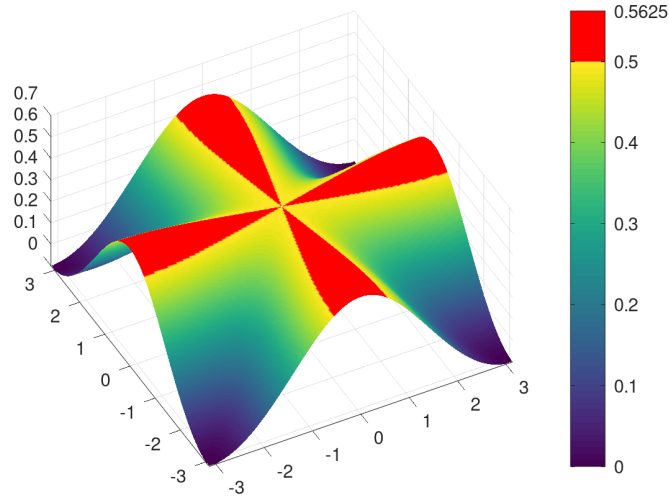


FIGURE 4.  $d = 2$ . Eigenvalue  $\lambda_+$  on  $Q$ .

In any case, we have proved that in dimension  $d = 2$  for any bounded open set  $\Omega \subset \mathbb{R}^2$  and any  $\lambda \in (-0.5471, -0.5) \cup (0.5, 0.5471)$  the delta-delta scheme  $(\lambda \mathbb{I} - T^N)U = F$  is unstable in  $\ell^2(\Omega^N; \mathbb{C}^2)$  as  $N \rightarrow \infty$ , whereas the integral equation  $(\lambda \mathbb{I} - A_\Omega)u = f$  is well posed in  $L^2(\Omega; \mathbb{C}^2)$ .

**3.5.4. Dimension  $d = 3$ .** The three eigenvalues  $\lambda_j$  of  $F(\tau)$  satisfy  $\lambda_1 + \lambda_2 + \lambda_3 = 0$ . Numerically, one sees that the minimal and maximal values are attained on the intersection of the boundary of  $Q = [-\pi, \pi]^3$  with the coordinate planes. In Figure 5 we show a graph of the three eigenvalues on the line  $\{(\pi, y, 0) \mid y \in [-\pi, \pi]\}$ . The values  $-\frac{1}{3}$  and  $\frac{2}{3}$  are shown as dashed lines.

This suggests  $\Lambda_-^{(3)} = \min \text{Sp}(F((\pi, \pi, 0)))$  and  $\Lambda_+^{(3)} = \max \text{Sp}(F((\pi, 0, 0)))$ .

The computed values are

$$(3.51) \quad \Lambda_-^{(3)} = -0.4260241507272727, \quad \Lambda_+^{(3)} = 0.77090222227747195.$$

This implies a length of  $W(T)$  of  $\Lambda_+^{(3)} - \Lambda_-^{(3)} = 1.1969263735019922$  instead of 1, which is the length of  $W(A)$ , thus an overshoot of almost 20%.

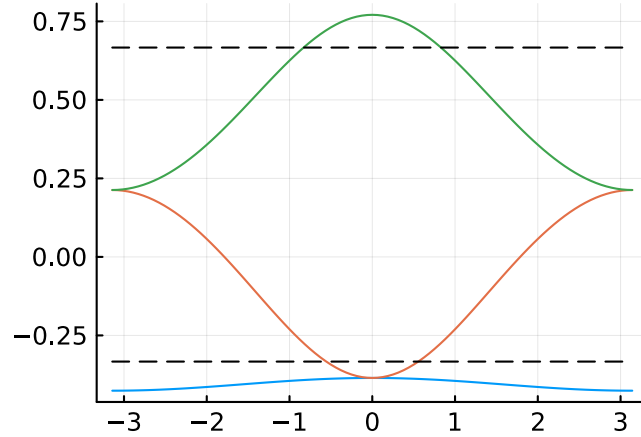


FIGURE 5.  $d = 3$ . Eigenvalues of  $F(\tau)$  on middle line of face of  $Q$ .

Under this assumption, one can write simple series expansions for the numbers  $\Lambda_{\pm}^{(3)}$ . If all the coordinates of  $\tau$  are 0 or  $\pi$ , then the off-diagonal elements of the matrix  $F(\tau)$  vanish and the 3 eigenvalues are the diagonal elements. Therefore the Fourier series for  $F(\tau)$  gives

$$(3.52) \quad \Lambda_{-}^{(3)} = \sum_{m \in \mathbb{Z}^3, m \neq 0} \frac{(-1)^{m_1+m_2}}{4\pi} \frac{m_1^2 + m_2^2 - 2m_3^2}{(m_1^2 + m_2^2 + m_3^2)^{\frac{5}{2}}}, \quad \Lambda_{+}^{(3)} = \sum_{m \in \mathbb{Z}^3, m \neq 0} \frac{(-1)^{m_3}}{4\pi} \frac{m_1^2 + m_2^2 - 2m_3^2}{(m_1^2 + m_2^2 + m_3^2)^{\frac{5}{2}}}.$$

These sums, although not absolutely convergent, appear to converge quite well in the sense of partial sums over cubes,

$$\sum_{m \in \mathbb{Z}^3, m \neq 0} = \lim_{N \rightarrow \infty} \sum_{\max_j |m_j| \leq N, m \neq 0}.$$

We do not know whether explicit expressions for these sums exist.

By means of the Clausius-Mossotti relation (1.8) one can express the stability results equivalently in terms of the relative permittivity  $\epsilon_r$ . Let

$$(3.53) \quad \epsilon_{\min} = \frac{3\Lambda_{+} - 2}{1 + 3\Lambda_{+}} = 0.0943961 \dots, \quad \epsilon_{\max} = \frac{3\Lambda_{-} - 2}{1 + 3\Lambda_{-}} = 11.788555 \dots.$$

The numerical range  $\lambda \in [-\frac{1}{3}, \frac{2}{3}]$  of the quasi-static Maxwell volume integral operator corresponds to  $\epsilon_r \leq 0$ . The volume integral equation is therefore well posed in  $L^2(\Omega)$  if the relative permittivity  $\epsilon_r$  is either non-real or positive.

On the other hand, the corresponding DDA scheme is stable in  $\ell^2(\mathbb{Z}^3)$  if and only if  $\epsilon_r$  is either non-real or contained in the interval  $(\epsilon_{\min}, \epsilon_{\max})$ . For  $\epsilon_r \in (0, \epsilon_{\min}] \cup [\epsilon_{\max}, \infty)$  the integral equation (and therefore the dielectric scattering problem) is well-posed, but the DDA scheme is unstable.

To conclude this discussion, we show in Table 4 the result of some computations for the spectrum of the system matrix  $T^N$  for a cube in three dimensions. One can see convergence to the expected values (3.51), even for rather modest values of  $N$ . Compare also [20, FIG. 8].



$N$	$\lambda_{\max}(T^N)$	$\lambda_{\min}(T^N)$	$\lambda_{\max}(T^N) - \lambda_{\min}(T^N)$
4	0.67730278666935	-0.3896455148525014	1.06694830152185
8	0.73653727456221	-0.4130173055963489	1.14955458015856
12	0.75323748914578	-0.4193953119966648	1.17263280114245
16	0.76017444184544	-0.4220149407429199	1.18218938258836

TABLE 4. Computation of  $\max(\text{Sp}(T^N))$  and  $\min(\text{Sp}(T^N))$ , Example 5

**Acknowledgment.** This work was partially supported by a grant from the *Niels Hendrik Abel Board*. The authors acknowledge support of the Centre Henri Lebesgue ANR-11-LABX-0020-01.

## REFERENCES

- [1] ABRAMOWITZ, M., AND STEGUN, I. A., Eds. *Handbook of mathematical functions with formulas, graphs and mathematical tables*. Washington: U.S. Department of Commerce. xiv, 1046 pp., 1964.
- [2] AMMARI, H., FITZPATRICK, B., KANG, H., RUIZ, M., YU, S., AND ZHANG, H. *Mathematical and computational methods in photonics and phononics*, vol. 235. Providence, RI: American Mathematical Society (AMS), 2018.
- [3] CHAUMET, P. C. The discrete dipole approximation: A review. *Mathematics* 10, 17 (2022).
- [4] COSTABEL, M., DARRIGRAND, E., AND SAKLY, H. The essential spectrum of the volume integral operator in electromagnetic scattering by a homogeneous body. *Comptes Rendus Mathématique* 350 (2012), 193–197.
- [5] COSTABEL, M., DARRIGRAND, E., AND SAKLY, H. Volume integral equations for electromagnetic scattering in two dimensions. *Comput. Math. Appl.* 70, 8 (2015), 2087–2101.
- [6] ESSMANN, U., PERERA, L., BERKOWITZ, M. L., DARDEN, T., LEE, H., AND PEDERSEN, L. G. A smooth particle mesh Ewald method. *J. Chem. Phys.* 103 (1995), 8577–8593.
- [7] EWALD, P. P. Die Berechnung optischer und elektrostatischer Gitterpotentiale. *Ann. der Phys. (4)* 64 (1921), 253–287.
- [8] GEL’FAND, I. M., AND SHILOV, G. E. *Generalized functions. Vol. I: Properties and operations*. Translated by E. Saletan. 1964.
- [9] JACKSON, J. D. *Classical Electrodynamics*, 3rd ed. John Wiley & Sons, Inc., 1999.
- [10] KIRSCH, A. An integral equation approach and the interior transmission problem for Maxwell’s equations. *Inverse Probl. Imaging* 1, 1 (2007), 159–179.
- [11] KOPPELMAN, W., AND PINCUS, J. D. Spectral representations for finite Hilbert transformations. *Math. Z.* 71 (1959), 399–407.
- [12] PÓLYA, G., AND SZEGÖ, G. *Isoperimetric inequalities in mathematical physics*, vol. 27 of *Ann. Math. Stud.* Princeton University Press, Princeton, NJ, 1951.
- [13] PURCELL, E. M., AND PENNYPACKER, C. R. Scattering and adsorption of light by nonspherical dielectric grains. *Astrophys. J.* 186 (1973), 705–714.
- [14] RAHOLA, J. On the eigenvalues of the volume integral operator of electromagnetic scattering. *SIAM J. Sci. Comput.* 21, 5 (2000), 1740–1754.
- [15] SMUNEV, D. A., CHAUMET, P. C., AND YURKIN, M. A. Rectangular dipoles in the discrete dipole approximation. *Journal of Quantitative Spectroscopy and Radiative Transfer* 156, 0 (2015), 67 – 79.
- [16] SÖHNGEN, H. Zur Theorie der endlichen Hilbert-Transformation. *Math. Z.* 60 (1954), 31–51.
- [17] TRICOMI, F. G. *Integral equations*. Pure and Applied Mathematics. Vol. V. Interscience Publishers, Inc., New York; Interscience Publishers Ltd., London, 1957.
- [18] YURKIN, M. A., AND HOEKSTRA, A. G. The discrete dipole approximation: An overview and recent developments. *J. Quant. Spectrosc. Radiat. Transf.* 106, 1 (2007), 558–589.

- [19] YURKIN, M. A., MALTSEV, V. P., AND HOEKSTRA, A. G. Convergence of the discrete dipole approximation. I. Theoretical analysis. *J. Opt. Soc. Am. A* 23, 10 (Oct 2006), 2578–2591.
- [20] YURKIN, M. A., MIN, M., AND HOEKSTRA, A. G. Application of the discrete dipole approximation to very large refractive indices: Filtered coupled dipoles revived. *Phys. Rev. E* 82 (Sep 2010), 036703.
- [21] ZUCKER, I. J. The summation of series of hyperbolic functions. *SIAM J. Math. Anal.* 10 (1979), 192–206.

UNIV. RENNES, CNRS, IRMAR - UMR 6625, F-35000 RENNES, FRANCE

*E-mail address:* Martin.Costabel@univ-rennes1.fr

*E-mail address:* Monique.Dauge@univ-rennes1.fr

DEPARTMENT OF MATHEMATICS, IASBS, GAVAZANG ROAD, ZANJAN, IRAN

*E-mail address:* knedaiasl85@gmail.com, nedaiasl@iasbs.ac.ir