

Preconditioning by inverting the Laplacian: an analysis of the eigenvalues

BJØRN FREDRIK NIELSEN[†], ASLAK TVEITO[‡] AND GANG ZHENG
*Simula Research Laboratory, PO Box 134, N-1325 Lysaker, Norway and
Department of Informatics, University of Oslo, PO Box 1080, N-0316*

AND

WOLFGANG HACKBUSCH[§]
*Max Planck Institute for Mathematics in the Sciences,
Inselstrasse 22–26, 04103 Leipzig, Germany*

[Received on 30 April 2008; revised on 29 August 2008]

We study the eigenvalues of the operator generated by using the inverse of the Laplacian as a preconditioner for self-adjoint second-order elliptic partial differential equations with smooth coefficients. It is well known that the spectral condition number of the preconditioned operator can be bounded by $\max k(x)/\min k(x)$, where k is the uniformly positive coefficient of the second-order elliptic equation. The purpose of this paper is to study the spectrum of the preconditioned operator. We will show that there is a strong relation between the spectrum of this operator and the range of the coefficient function. In the continuous case we prove, both for mappings defined on Sobolev spaces and in terms of generalized functions, that the spectrum of the preconditioned operator contains the range of the coefficient function k . In the discrete case we indicate by numerical examples that the entire discrete spectrum is approximately given by values of k .

Keywords: elliptic differential equations; eigenvalues; preconditioning; continuous coefficients.

1. Introduction

Consider a self-adjoint second-order elliptic partial differential equation (PDE) of the form

$$-\nabla \cdot (k(x, y) \nabla v) = f \quad \text{for } (x, y) \in \Omega, \quad (1.1)$$

equipped with suitable boundary conditions. Here k is a uniformly positive and smooth function and f is a given function in $L^2(\Omega)$. It is well known that a straightforward discretization using the finite-element method leads to a discrete system of the form

$$A_h v_h = f_h, \quad (1.2)$$

where h denotes the mesh parameter. In a similar manner, we can discretize (1.1) in the case of $k \equiv 1$, corresponding to the Poisson equation, and obtain a linear system of the form

$$L_h w_h = f_h. \quad (1.3)$$

[†]Corresponding author. Email: bjornn@simula.no

[‡]Email: aslak@simula.no

[§]Email: wh@mis.mpg.de

It is also known that, if L_h^{-1} is used as a preconditioner for A_h in the process of solving (1.2) numerically, then the condition number of the preconditioned operator, formally denoted by $L_h^{-1}A_h$, is bounded by

$$\kappa \leq \frac{\sup_{(x,y) \in \Omega} k(x,y)}{\inf_{(x,y) \in \Omega} k(x,y)}.$$

Furthermore, it is known that the number of conjugated gradient (CG) iterations needed to solve the problem is $\mathcal{O}(\sqrt{\kappa})$, and thus the number of iterations is independent of the mesh parameter h (see, e.g., Axelsson, 1994).

Even though the spectral condition number κ of $L_h^{-1}A_h$ is bounded independently of h , many CG iterations may still be needed if the variations in k are large. For coefficients with jump discontinuities this problem has partly been successfully solved by introducing various preconditioners. In such cases the methods and analyses are commonly based on some sort of clustering effect of the eigenvalues of the preconditioned operator (see Cai *et al.*, 1999; Graham & Hagger, 1999).

The aim of the present paper is to study this problem for equations involving continuous coefficients with large variations. We will do this by providing insight into the spectrum of the preconditioned operator $L_h^{-1}A_h$.

Our main results can roughly be described as follows. Let $\lambda_h = (\lambda_1, \lambda_2, \dots)$ denote a vector containing all of the eigenvalues of $L_h^{-1}A_h$ sorted in a non-decreasing manner. These eigenvalues can be computed by solving the generalized eigenvalue problem

$$A_h u_h = \lambda L_h u_h. \quad (1.4)$$

Similarly, we let $\mu_h = (\mu_1, \mu_2, \dots)$ denote a vector containing the values of the coefficient function k evaluated at the mesh points and sorted in a non-decreasing manner. We will present numerical experiments indicating that

$$\max_j |\lambda_j - \mu_j| = \mathcal{O}(h).$$

In the continuous case the generalized eigenvalue problem is to find a non-trivial eigenfunction u and an eigenvalue λ such that

$$\nabla \cdot (k(x,y) \nabla u) = \lambda \nabla^2 u \quad \text{for } (x,y) \in \Omega \quad (1.5)$$

for a suitable set of boundary conditions. We will prove rigorously, both for the associated operators defined on Sobolev spaces and in terms of distribution theory, that all of the values of the coefficient function k are contained in the spectrum of the preconditioned operator $L^{-1}A$.

The implication of this is that we obtain a deeper understanding of the convergence properties of the CG method since the convergence is completely determined by the spectrum of the preconditioned operator (see, e.g., Stoer & Bulirsch, 1993). More precisely, if the coefficient k is continuous and has large variations then no clustering effect of the eigenvalues of $L_h^{-1}A_h$ will occur—indicating that efficient methods for such problems must be based on some other kind of property of the involved equation. This means that optimal preconditioners are well suited to obtaining h -independent bounds on the number of CG iterations, but such methods will still generate a large number of CG iterations for continuous coefficients k that have large variations. As mentioned above, the case of a piecewise constant k with large variations has been successfully solved.

For a general introduction to preconditioners and iterative methods we refer to Axelsson (1994), Chan & Mathew (1994) and Hackbusch (1994). The use of fast solvers for the Laplace equation as preconditioners for the CG method has been used for many years. For a thorough discussion on how the

eigenvalues of the operator influence the convergence properties of the CG method we refer to Axelsson & Lindskog (1986a,b). Numerical schemes for equations of the form (1.1) with large jumps in the coefficients were studied by, for example, Bramble *et al.* (1991), Dryja (1994), Dryja *et al.* (1996) and Bebendorf & Hackbusch (2003).

The convergence properties of the CG method are also heavily influenced by both the right-hand side f of (1.1) and the start vector of the iteration. These issues were thoroughly analysed by several scientists (see Naiman *et al.* (1997), Cai *et al.* (1999), Naiman & Engelberg (2000) and Beckermann & Kuijlaars (2001, 2002) for further information).

To shed some light onto the generalized eigenvalue problem (1.5), a series of numerical experiments is presented in Section 2. In Section 3 we study the continuous eigenvalue problem for operators defined on appropriate Sobolev spaces, and prove that, for any given point $(x_0, y_0) \in \Omega$, the value $k(x_0, y_0)$ of the coefficient function at this point is an eigenvalue for the preconditioned operator, provided that k is continuous at (x_0, y_0) . Furthermore, Section 4 contains an analysis of this problem in terms of distribution theory, where also explicit formulas for the eigenfunctions are presented. Some conclusions are discussed in Section 5.

2. Numerical experiments

The purpose of this section is to present some numerical experiments that indicate the properties of the preconditioned operator $L_h^{-1}A_h$. These examples provide motivation for the rigorous analysis presented in Sections 3 and 4 below.

Consider the equation

$$-\nabla \cdot (k(x, y)\nabla v) = f \quad \text{for } (x, y) \in \Omega = (0, 1) \times (0, 1), \quad (2.1)$$

with boundary condition

$$\frac{\partial v}{\partial n} = 0 \quad \text{at } \partial\Omega,$$

where n denotes a normal vector for the boundary of Ω , and we assume that f satisfies the solvability condition

$$\int_{\Omega} f \, dx = 0$$

(see, for example, Marti (1986) for further details).

We discretize this problem in a standard manner using linear elements on a uniform mesh (see, e.g., Langtangen, 1999). The discretization leads to a linear system of the form

$$A_h v_h = f_h. \quad (2.2)$$

We also discretize the problem in the case of $k \equiv 1$, and obtain a linear system of the form

$$L_h w_h = f_h. \quad (2.3)$$

Our aim is to compute the eigenvalues of the preconditioned operator $L_h^{-1}A_h$, i.e., we want to find the eigenvalues of the generalized eigenvalue problem

$$A_h u_h = \lambda L_h u_h. \quad (2.4)$$

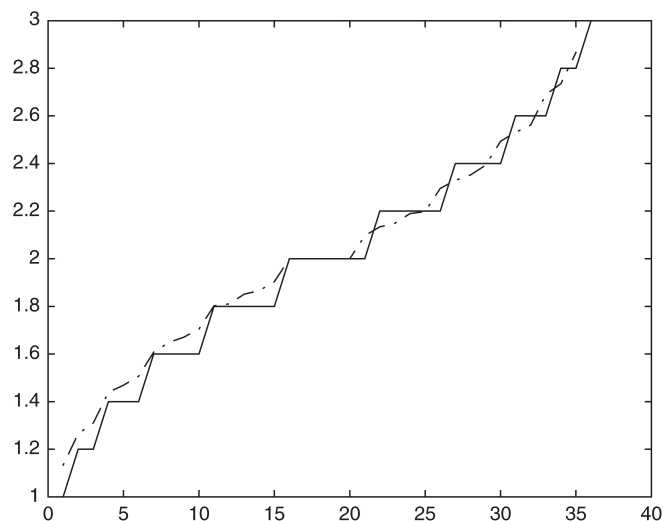


FIG. 1. The sorted array of eigenvalues (dashed-dotted line) and of the coefficient function k evaluated at the grid points (solid line). In this case $k = k_1$ and the mesh size is $h = 0.2$.

This is done by generating the matrices A_h and L_h using Diffpack (see Langtangen, 1999). The matrices¹ are stored and entered into Matlab, which is used to solve the problem (2.4).

In our numerical experiments we use three versions of the coefficient function k :

$$\begin{aligned} k_1(x,y) &= 1 + x + y, \\ k_2(x,y) &= 2 + \sin(2\pi e^x \cos(yx)), \\ k_3(x,y) &= 1 + 50e^{-50[(x-0.5)^2 + (y-0.5)^2]}. \end{aligned}$$

As above, we let $\lambda_h = (\lambda_1, \lambda_2, \dots)$ denote a vector containing all of the eigenvalues of $L_h^{-1}A_h$, i.e., of (2.4), sorted in a non-decreasing manner. Similarly, we let $\mu_h = (\mu_1, \mu_2, \dots)$ denote a vector containing the values of the coefficient function k evaluated at the mesh points and sorted in a non-decreasing manner.

In Figs 1 and 2 we plot the eigenvalues λ_h and the function values μ_h in the case of $k = k_1$ using $h = 0.2$ and $h = 0.05$, i.e., we have $n = 36$ and $n = 441$ unknowns. We note that the graphs are quite similar. Similar plots for $k = k_2$ and $k = k_3$ are presented in Figs 3–6.

¹It should be noted that, due to the boundary conditions, both matrices are singular. If a solution to the boundary-value problem is to be computed then it is common to add the following additional constraint:

$$\int_{\Omega} v dx = 0.$$

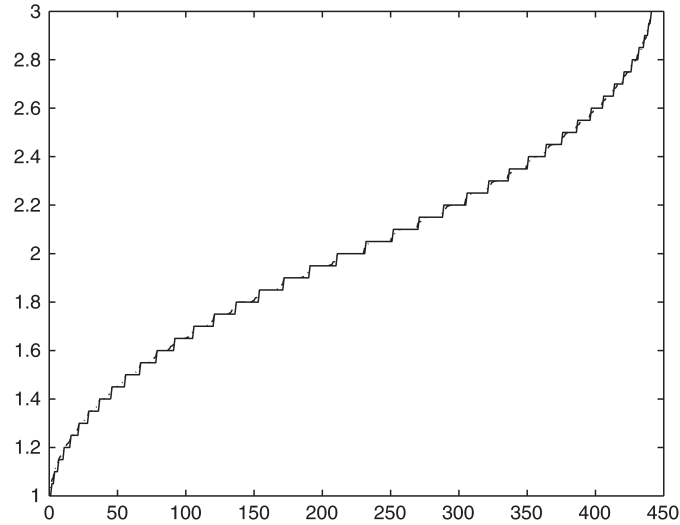


FIG. 2. The sorted array of eigenvalues (dashed-dotted line) and of the coefficient function k evaluated at the grid points (solid line). In this case $k = k_1$ and the mesh size is $h = 0.05$.

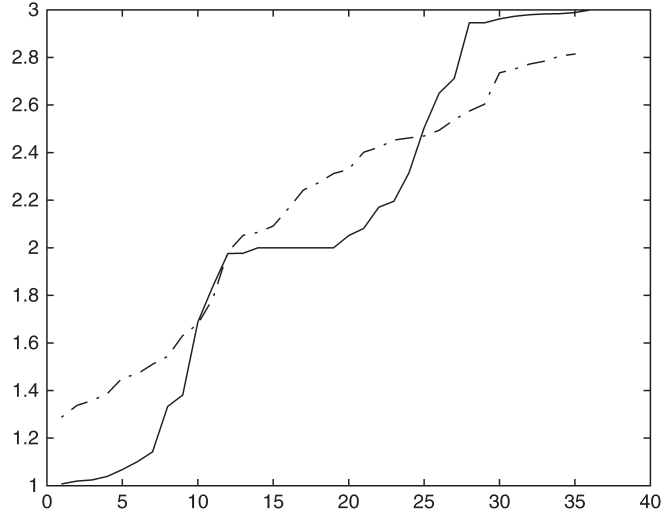


FIG. 3. The sorted array of eigenvalues (dashed-dotted line) and of the coefficient function k evaluated at the grid points (solid line). In this case $k = k_2$ and the mesh size is $h = 0.2$.

The numerical results are summarized in Tables 1–3, where we show the behaviour of $\max_j |\lambda_j - \mu_j|$ for various values of h . Clearly, these tables indicate that μ_h is an $\mathcal{O}(h)$ -approximation of λ_h . The rate of convergence is computed by comparing the results obtained for two successive values of h and by assuming that the difference $\max_j |\lambda_j - \mu_j|$ has the form ch^α , where c is a constant and α is the rate. (Ideally, for $k = k_3$ further experiments on finer meshes, i.e., with $h < 0.0125$, should be performed.

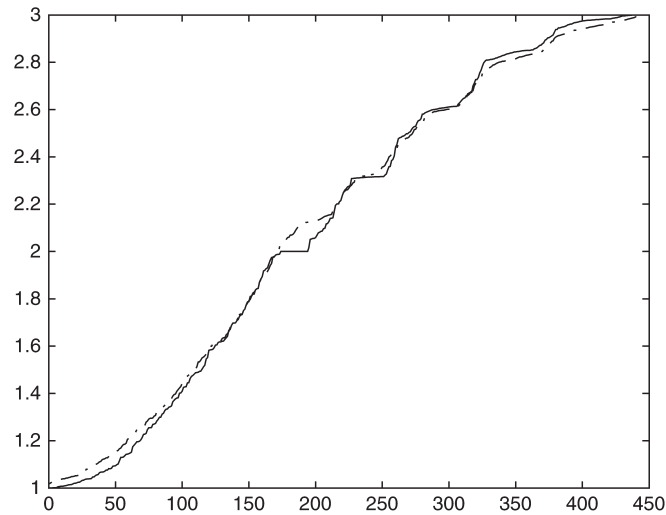


FIG. 4. The sorted array of eigenvalues (dashed–dotted line) and of the coefficient function k evaluated at the grid points (solid line). In this case $k = k_2$ and the mesh size is $h = 0.05$.

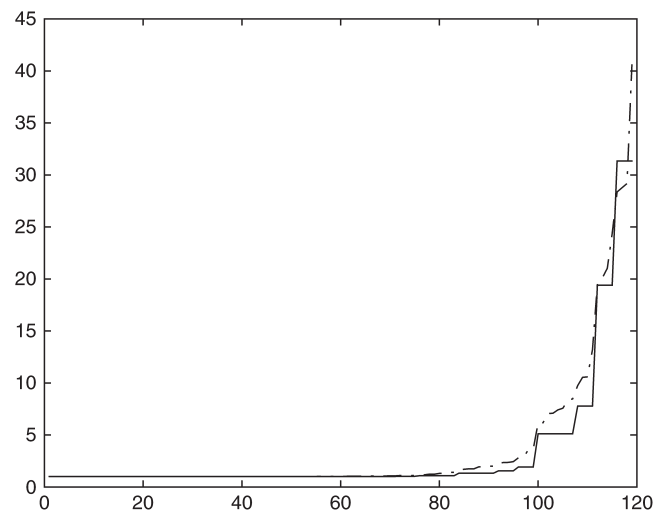


FIG. 5. The sorted array of eigenvalues (dashed–dotted line) and of the coefficient function k evaluated at the grid points (solid line). In this case $k = k_3$ and the mesh size is $h = 0.1$.

However, the Matlab algorithms for computing the eigenvalues of $L_h^{-1}A_h$ are very CPU and memory demanding, and such investigations were therefore not possible, within reasonable time limits, on our computers.

We have not succeeded in showing rigorously that μ_h defines an $\mathcal{O}(h)$ -approximation of λ_h —this is, as far as we know, still an open problem. However, in the continuous case, for operators defined on appropriate Sobolev spaces, we have proven that $k(x_0, y_0)$ is an eigenvalue for the preconditioned operator, provided that k is continuous at (x_0, y_0) . This issue is treated in detail in Section 3.

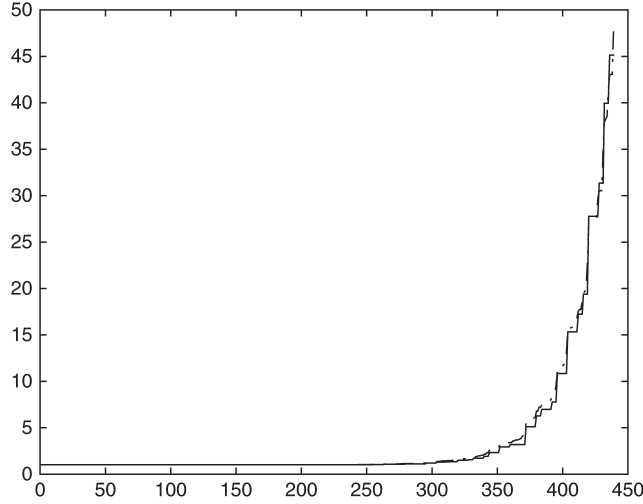


FIG. 6. The sorted array of eigenvalues (dashed–dotted line) and of the coefficient function k evaluated at the grid points (solid line). In this case $k = k_3$ and the mesh size is $h = 0.05$.

TABLE 1. *The numerical results obtained for the generalized eigenvalue problem (2.4) with $k(x, y) = k_1(x, y) = 1 + x + y$. In this table, as well as in Tables 2 and 3, λ_j , μ_j , h and n represent the eigenvalue, the value of the coefficient function evaluated at the mesh point, the mesh width and the number of unknowns in the discrete eigenvalue problem, respectively*

h	n	$\max_j \lambda_j - \mu_j $	$\max_j \lambda_j - \mu_j /h$	Rate
0.2	36	0.1333	0.6667	–
0.1	121	0.0667	0.6667	0.9989
0.05	441	0.0333	0.6667	1.0022
0.025	1681	0.0167	0.6667	0.9957
0.0125	6561	0.0083	0.6667	1.0087

TABLE 2. *The numerical results obtained for the generalized eigenvalue problem (2.4) with $k(x, y) = k_2(x, y) = 2 + \sin(2\pi e^x \cos(yx))$*

h	n	$\max_j \lambda_j - \mu_j $	$\max_j \lambda_j - \mu_j /h$	Rate
0.2	36	0.3847	1.9235	–
0.1	121	0.2009	2.0090	0.9373
0.05	441	0.1234	2.4672	0.7031
0.025	1681	0.0633	2.5325	0.9631
0.0125	6561	0.0317	2.5363	0.9977

TABLE 3. *The numerical results obtained for the generalized eigenvalue problem (2.4) with $k(x, y) = k_3(x, y) = 1 + 50e^{-50[(x-0.5)^2 + (y-0.5)^2]}$*

h	n	$\max_j \lambda_j - \mu_j $	$\max_j \lambda_j - \mu_j /h$	Rate
0.2	36	5.5235	27.6175	–
0.1	121	5.4286	54.2855	0.025
0.05	441	3.6823	73.6464	0.56
0.025	1681	2.3656	94.6229	0.6384
0.0167	3721	1.1833	71.0001	1.7169
0.0125	6561	0.8723	69.7820	1.0526

3. Operators defined on Sobolev spaces

In this section we will study some properties of a generalized eigenvalue problem of the following form: find a number λ and a function u such that

$$\nabla \cdot (k \nabla u) = \lambda \Delta u \quad \text{in } \Omega. \quad (3.1)$$

Our aim is to analyse this equation in terms of operators defined on Sobolev spaces.

To this end, let us consider a prototypical elliptic boundary-value problem of the form

$$\begin{aligned} -\nabla \cdot (k \nabla v) &= f \quad \text{in } \Omega, \\ v &= 0 \quad \text{on } \partial\Omega, \end{aligned} \quad (3.2)$$

where Ω , with boundary $\partial\Omega$, is some open Lipschitz domain contained in a Euclidean space \mathbb{R}^n . In addition, k is assumed to be a uniformly positive and bounded function defined on Ω , i.e.,

$$k \in L^\infty(\Omega), \quad (3.3)$$

$$0 < b \leq k(x) \leq B \quad \text{for all } x \in \Omega, \quad (3.4)$$

where b and B are given positive constants.

Throughout this section we will, for the sake of easy notation, consider elliptic PDEs with homogeneous Dirichlet boundary conditions (cf. (3.2)). However, it is not difficult to modify the arguments presented below to also cover cases involving Neumann conditions.

3.1 Notation

Let $H_0^1(\Omega)$, with inner product $(\cdot, \cdot)_1$ and norm $\|\cdot\|_{H^1(\Omega)}$, denote the classical Sobolev space of functions defined on Ω with zero trace at $\partial\Omega$. According to the Riesz representation theorem, there exist linear operators $A, L \in \mathcal{L}(H_0^1(\Omega))$ such that

$$(A\varphi, \psi)_1 = \int_{\Omega} \nabla \psi \cdot (k \nabla \varphi) \, dx \quad \text{for all } \psi, \varphi \in H_0^1(\Omega), \quad (3.5)$$

$$(L\varphi, \psi)_1 = \int_{\Omega} \nabla \psi \cdot \nabla \varphi \, dx \quad \text{for all } \psi, \varphi \in H_0^1(\Omega). \quad (3.6)$$

With this notation at hand, the weak form of (3.2) may be written in the following form: find $v \in H_0^1(\Omega)$ such that

$$(Av, \psi)_1 = \int_{\Omega} f \psi \, dx \quad \text{for all } \psi \in H_0^1(\Omega).$$

Furthermore, L represents the ('weak') Laplacian defined on $H_0^1(\Omega)$.

3.2 The analysis

Inspired by our findings in Section 2, we will now analyse the spectrum²

$$\text{sp}(L^{-1}A) = \{\lambda \in \mathbf{C}; (\lambda I - L^{-1}A) \text{ is not invertible}\}$$

of the operator

$$L^{-1}A : H_0^1(\Omega) \rightarrow H_0^1(\Omega).$$

For the sake of convenience, let K denote the set of points in Ω at which the coefficient function k is continuous, i.e.,

$$K = \{x \in \Omega; k \text{ is continuous at } x\}. \quad (3.7)$$

Our main result can now be formulated as follows.

THEOREM 3.1 Let A and L be the operators defined in (3.5) and (3.6).

(a) If (3.3) and (3.4) hold then

$$k(x) \in \text{sp}(L^{-1}A) \quad \text{for all } x \in K,$$

where K is the set of points defined in (3.7).

(b) In particular, if (3.4) holds and k is continuous throughout Ω , i.e., $k \in C(\Omega)$, then

$$k(x) \in \text{sp}(L^{-1}A) \quad \text{for all } x \in \Omega.$$

Proof. Let $\tilde{x} \in K$ be arbitrary and assume that \tilde{x} is such that $\tilde{\lambda} = k(\tilde{x}) \notin \text{sp}(L^{-1}A)$.

Clearly, there exists a set of functions $\{v_r\}_{r \in \mathbb{R}_+}$ satisfying³

$$\text{supp}(v_r) \subset \tilde{x} + U_r, \quad (3.8)$$

$$\|v_r\|_{H^1(\Omega)} = 1, \quad (3.9)$$

where

$$U_r = \{\mathbf{z} \in \mathbb{R}^n; |\mathbf{z}| \leq r\}.$$

Let

$$u_r = (\tilde{\lambda}I - L^{-1}A)v_r \quad \text{for } r \in \mathbb{R}_+. \quad (3.10)$$

²Here $I : H_0^1(\Omega) \rightarrow H_0^1(\Omega)$ denotes the identity operator, i.e.,

$$I\psi = \psi \quad \text{for all } \psi \in H_0^1(\Omega).$$

³Note that no limit of v_r as $r \rightarrow 0$ is needed in this proof. Only the existence of a set of functions satisfying (3.8) and (3.9) is required.

Since $\tilde{\lambda} \notin \text{sp}(L^{-1}A)$, it follows that $(\tilde{\lambda}I - L^{-1}A)$ is invertible, and we find that

$$\|v_r\|_{H^1(\Omega)} = \|(\tilde{\lambda}I - L^{-1}A)^{-1}u_r\|_{H^1(\Omega)} \leq \|(\tilde{\lambda}I - L^{-1}A)^{-1}\| \|u_r\|_{H^1(\Omega)}. \quad (3.11)$$

By (3.10), we have

$$u_r = \tilde{\lambda}Iv_r - L^{-1}Av_r$$

or

$$Lu_r = \tilde{\lambda}Lv_r - Av_r.$$

Hence it follows that

$$(Lu_r, u_r)_1 = \tilde{\lambda}(Lv_r, u_r)_1 - (Av_r, u_r)_1,$$

and then, from the definition of L and A and by the fact that $\text{supp}(v_r) \subset \tilde{x} + U_r$, we find that

$$\int_{\Omega} \nabla u_r \cdot \nabla u_r \, dx = \tilde{\lambda} \int_{\tilde{x}+U_r} \nabla u_r \cdot \nabla v_r \, dx - \int_{\tilde{x}+U_r} \nabla u_r \cdot (k \nabla v_r) \, dx$$

or

$$\int_{\Omega} |\nabla u_r|^2 \, dx = \int_{\tilde{x}+U_r} \nabla u_r \cdot ([\tilde{\lambda} - k] \nabla v_r) \, dx.$$

Next, by the Cauchy–Schwarz inequality, we have

$$\int_{\Omega} |\nabla u_r|^2 \, dx \leq \left(\int_{\Omega} |\nabla u_r|^2 \, dx \right)^{1/2} \left(\int_{\tilde{x}+U_r} (\tilde{\lambda} - k)^2 |\nabla v_r|^2 \, dx \right)^{1/2},$$

and, consequently,

$$\left(\int_{\Omega} |\nabla u_r|^2 \, dx \right)^{1/2} \leq \text{ess sup}_{x \in \tilde{x}+U_r} |\tilde{\lambda} - k(x)| \|v_r\|_{H^1(\Omega)} = \text{ess sup}_{x \in \tilde{x}+U_r} |k(\tilde{x}) - k(x)|, \quad (3.12)$$

where the last equality follows from (3.9). Since k is continuous at \tilde{x} , it follows that

$$\lim_{r \rightarrow 0} \text{ess sup}_{x \in \tilde{x}+U_r} |k(\tilde{x}) - k(x)| = 0.$$

From (3.12) and Poincaré's inequality we thus conclude that there exists a constant $r^* \in \mathbb{R}_+$ such that

$$\|u_r\|_{H^1(\Omega)} < \frac{1}{2 \|(\tilde{\lambda}I - L^{-1}A)^{-1}\|} \quad \text{for all } r \in (0, r^*). \quad (3.13)$$

Finally, (3.11) and (3.13) imply that

$$\|v_r\|_{H^1(\Omega)} < \frac{1}{2} \quad \text{for all } r \in (0, r^*),$$

which is a contradiction to (3.9). Hence we conclude that $k(\tilde{x}) \in \text{sp}(L^{-1}A)$.

This completes the proof of part (a) of the theorem. Part (b) is a trivial consequence of part (a). \square

If k is continuous then we thus conclude that the range of k is indeed contained in the spectrum of $L^{-1}A$, which is in agreement with the results of our numerical experiments. Moreover, for problems involving discontinuous coefficient functions, i.e., $k \notin C(\Omega)$, we can still conclude that $k(x)$ is an eigenvalue for the preconditioned operator $L^{-1}A$ at every point x at which k is continuous.

As mentioned in Section 1, for elliptic equations with coefficients with large jump discontinuities, high-quality preconditioners can sometimes be constructed due to some clustering effect of the eigenvalues. In the case of continuous coefficients our investigations indicate that such an effect is not likely to occur. In particular, if the inverse Laplacian is used as a preconditioner then the eigenvalues will not cluster.

The proof of Theorem 3.1 is not of a constructive nature. Formulas for the eigenfunctions are not provided. To further shed some light onto the generalized eigenvalue problem (3.1) we will now consider it from a distributional point of view.

4. Generalized eigenfunctions and eigenvalues

As mentioned above, our aim is to study the eigenvalue problem

$$\nabla \cdot (k \nabla u) = \lambda \Delta u \tag{4.1}$$

in terms of distribution theory. More precisely, we will not only prove that $\lambda = k(x, y)$, for all $(x, y) \in \Omega$, are eigenvalues but also present explicit formulas for the associated generalized eigenfunctions. As in Section 3, we will assume that Ω is an open Lipschitz domain.

4.1 Preliminaries and notation

In our analysis of this eigenvalue problem the classical mollifier functions (see, e.g., Marti, 1986) will play an important role. It is defined by a non-negative and symmetric function $\omega \in C^\infty(\mathbb{R})$ satisfying⁴

$$\int_{\mathbb{R}} \omega(x) dx = 1,$$

$$\omega(x) \equiv 0 \quad \text{for } |x| \geq 1.$$

A family $\{\omega_\varepsilon\}_{\varepsilon \in \mathbb{R}_+}$ of mollifier functions are now defined by setting

$$\omega_\varepsilon = \frac{1}{\varepsilon} \omega\left(\frac{x}{\varepsilon}\right).$$

Clearly, these functions possess the following properties:

$$\omega_\varepsilon \in C^\infty, \quad \omega_\varepsilon \geq 0, \tag{4.2}$$

$$\omega_\varepsilon(x) = \omega_\varepsilon(-x), \tag{4.3}$$

$$\int_{\mathbb{R}} \omega_\varepsilon(x) dx = 1, \tag{4.4}$$

$$\omega_\varepsilon(x) \equiv 0 \quad \text{for } |x| \geq \varepsilon, \tag{4.5}$$

$$\omega_\varepsilon(x) \leq \frac{M}{\varepsilon} \quad \text{and} \quad |\omega'_\varepsilon(x)| \leq \frac{M}{\varepsilon^2}, \tag{4.6}$$

⁴The standard example of such a function is

$$\omega(x) = \begin{cases} c \exp((x^2 - 1)^{-1}) & \text{for } x \in (-1, 1), \\ 0 & \text{for } |x| \geq 1, \end{cases}$$

where $c^{-1} = \int_{-1}^1 \exp((x^2 - 1)^{-1}) dx$.

where M is a positive constant that is independent of ε . Next we define a family $\{H^\varepsilon\}_{\varepsilon \in \mathbb{R}_+}$ of approximate Heaviside functions

$$H^\varepsilon(x) = \int_{-\infty}^x \omega_\varepsilon(y) dy.$$

Note that

$$0 \leq H^\varepsilon(x) \leq 1 \quad \text{for all } x \in \mathbb{R}, \quad (4.7)$$

$$H^\varepsilon(x) \equiv 0 \quad \text{for } x \leq -\varepsilon, \quad H^\varepsilon(x) \equiv 1 \quad \text{for } x \geq \varepsilon, \quad (4.8)$$

$$(H^\varepsilon)'(x) = \omega_\varepsilon(x). \quad (4.9)$$

We have not been able to characterize the (generalized) eigenfunctions and eigenvalues satisfying (4.1) for all smooth coefficient functions k . However, we have been able to do so for a fairly large class of coefficient functions that we now describe. To this end, we define the following family \mathcal{Q} of smooth and uniformly positive functions defined on Ω :

$$\mathcal{Q} = \{k \in C^\infty(\overline{\Omega}); \exists m \in \mathbb{R}_+ \text{ such that } m \leq k(x, y) \text{ for all } (x, y) \in \overline{\Omega}\}.$$

It turns out that the generalized eigenfunctions satisfying (4.1) are characterized by the regions in the domain Ω on which the coefficient function k is constant, i.e., by the contour curves of k . Therefore, for each $k \in \mathcal{Q}$ and $(x_0, y_0) \in \Omega$, we introduce the set

$$S(k, (x_0, y_0)) = \{(x, y) \in \Omega; k(x, y) = k(x_0, y_0)\}. \quad (4.10)$$

With this notation at hand, we are ready to define the set \mathcal{K} of coefficient functions for which we are able to provide a detailed mathematical analysis of the eigenvalue problem (4.1) as follows:

$$\begin{aligned} \mathcal{K} = \{k \in \mathcal{Q}; \text{ for all } (x, y) \in \Omega, |S(k, (x, y))| = 0 \text{ or } S(k, (x, y)) \\ \text{contains at least one open and connected subset } G \text{ with } |G| > 0\}. \end{aligned} \quad (4.11)$$

Here $|S(k, (x, y))|$ and $|G|$ denote the measures of the respective sets. Roughly speaking, \mathcal{K} consists of those smooth functions that are ‘well behaved’ in a measure-theoretical sense. The need for these assumptions on the coefficient function will become apparent in the analysis below.

Clearly, the distributional form of the eigenvalue problem (4.1) can be written in the following form: find a number λ and a (possibly generalized) function u such that

$$\int_{\Omega} k \nabla u \cdot \nabla \phi \, d\Omega = \lambda \int_{\Omega} \nabla u \cdot \nabla \phi \, d\Omega \quad \text{for all } \phi \in C_0^\infty(\Omega), \quad (4.12)$$

where $C_0^\infty(\Omega)$ denotes the set of test functions, i.e., the set⁵ of smooth functions with compact support in Ω . This means that

$$\int_{\Omega} (k - \lambda)[u_x \phi_x + u_y \phi_y] \, d\Omega = 0 \quad \text{for all } \phi \in C_0^\infty(\Omega). \quad (4.13)$$

In the analysis below we will use the form (4.13) of the generalized eigenvalue problem. The analysis is divided, by certain properties of the coefficient function k , into three different cases.

⁵ $C_0^\infty(\Omega) = \{\psi \in C^\infty(\Omega); \exists \text{ a compact set } K \subset \Omega \text{ such that } \{x; \psi(x) \neq 0\} \subset K\}.$

4.2 The analysis

4.2.1 *Case I.* Let $k \in K$ and assume that $(x_0, y_0) \in \Omega$ is a point such that

$$k_x(x_0, y_0) \neq 0 \quad \text{or} \quad k_y(x_0, y_0) \neq 0, \quad (4.14)$$

$$|S(k, (x_0, y_0))| = 0. \quad (4.15)$$

We will now study the sequence of functions

$$H^\varepsilon(k(x, y) - k_0)$$

for $\varepsilon > 0$, where $k_0 = k(x_0, y_0)$. More precisely, we will show that k_0 is an eigenvalue with associated eigenfunction $H(k(x, y) - k_0)$, where H denotes the Heaviside function

$$H(z) = \begin{cases} 0 & \text{for } z \leq 0, \\ 1 & \text{for } z > 0. \end{cases} \quad (4.16)$$

Motivated by the discussion of the numerical experiments above and the form (4.13) of the generalized eigenvalue problem, we will, for an arbitrary test function $\phi \in C_0^\infty(\Omega)$, study the integral

$$\begin{aligned} I_\varepsilon &= \int_{\Omega} (k - k_0) [H_x^\varepsilon(k - k_0)\phi_x + H_y^\varepsilon(k - k_0)\phi_y] \, d\Omega \\ &= \int_{\Omega} (k - k_0) [k_x \omega_\varepsilon(k - k_0)\phi_x + k_y \omega_\varepsilon(k - k_0)\phi_y] \, d\Omega. \end{aligned} \quad (4.17)$$

If we apply the property (4.5) of the mollifier and define

$$S_\varepsilon = \{(x, y) \in \Omega; |k(x, y) - k(x_0, y_0)| < \varepsilon\} \quad (4.18)$$

then we find that

$$I_\varepsilon = \int_{S_\varepsilon} (k - k_0) [k_x \omega_\varepsilon(k - k_0)\phi_x + k_y \omega_\varepsilon(k - k_0)\phi_y] \, d\Omega.$$

Next, since $k_x, k_y, \phi_x, \phi_y \in L_\infty(\Omega)$, property (4.6) of the mollifier function ω_ε and the triangle inequality imply the existence of a positive constant c_1 , independent of ε , such that

$$\begin{aligned} |I_\varepsilon| &\leq \int_{S_\varepsilon} |k - k_0| [|k_x \omega_\varepsilon(k - k_0)\phi_x| + |k_y \omega_\varepsilon(k - k_0)\phi_y|] \, d\Omega \\ &\leq \int_{S_\varepsilon} |k - k_0| [\|k_x\|_\infty \|\phi_x\|_\infty M\varepsilon^{-1} + \|k_y\|_\infty \|\phi_y\|_\infty M\varepsilon^{-1}] \, d\Omega \\ &\leq c_1 \int_{S_\varepsilon} |k - k_0| \varepsilon^{-1} \, d\Omega. \end{aligned}$$

However, on the set S_ε (see (4.18)) we have $|k - k_0| < \varepsilon$ and we therefore conclude that

$$|I_\varepsilon| \leq c_1 \int_{S_\varepsilon} 1 \, d\Omega = c_1 |S_\varepsilon|, \quad (4.19)$$

where $|S_\varepsilon|$ denotes the measure of S_ε . Recall the definitions (4.10) and (4.18) of $S(k, (x_0, y_0))$ and S_ε , respectively. Clearly,

$$\bigcap_{\varepsilon > 0} S_\varepsilon = S(k, (x_0, y_0)),$$

and it is therefore natural to ask if the measure of S_ε converges toward the measure of $S(k, (x_0, y_0))$ as $\varepsilon \rightarrow 0$. This question is treated in detail in Appendix A and the answer to it is affirmative (see Lemma A.1). Thus assumption (4.15) implies that

$$\lim_{\varepsilon \rightarrow 0} |S_\varepsilon| = 0,$$

and we conclude that

$$\lim_{\varepsilon \rightarrow 0} |I_\varepsilon| = 0. \quad (4.20)$$

Having established that the integral defined in (4.17) tends toward zero as $\varepsilon \rightarrow 0$, we must now check whether or not the sequence of functions $\{H^\varepsilon(k - k_0)\}_{\varepsilon \in \mathbb{R}_+}$ has a well-defined, in the distributional sense, limit as $\varepsilon \rightarrow 0$. More precisely, we will show, as expected, that $H(k - k_0)$ is the limit. To this end, consider for an arbitrary test function $\phi \in C_0^\infty$ the integral

$$\begin{aligned} D_\varepsilon &= \int_{\Omega} H^\varepsilon(k - k_0) \phi \, d\Omega - \int_{\Omega} H(k - k_0) \phi \, d\Omega \\ &= \int_{\Omega} [H^\varepsilon(k - k_0) - H(k - k_0)] \phi \, d\Omega. \end{aligned}$$

From the property (4.8) of the approximate Heaviside function H^ε we find that

$$H^\varepsilon(k - k_0) = H(k - k_0) \quad \text{for all } (x, y) \text{ such that } |k(x, y) - k_0| \geq \varepsilon.$$

Hence

$$D_\varepsilon = \int_{S_\varepsilon} [H^\varepsilon(k - k_0) - H(k - k_0)] \phi \, d\Omega,$$

and then the property (4.7) and the definition (4.16) of the Heaviside function imply that

$$|D_\varepsilon| \leq \|\phi\|_\infty \int_{S_\varepsilon} d\Omega = \|\phi\|_\infty |S_\varepsilon|.$$

As discussed above, $|S_\varepsilon| \rightarrow 0$ as $\varepsilon \rightarrow 0$, and we conclude that

$$\lim_{\varepsilon \rightarrow 0} H^\varepsilon(k - k_0) = H(k - k_0)$$

in the distributional sense. From standard theory (see Griffel, 1981), for the derivatives of distributions, it follows that

$$\lim_{\varepsilon \rightarrow 0} H_x^\varepsilon(k - k_0) = H_x(k - k_0), \quad \lim_{\varepsilon \rightarrow 0} H_y^\varepsilon(k - k_0) = H_y(k - k_0).$$

Finally, by combining these convergence properties of the approximate Heaviside functions and (4.17) and (4.20) we find that

$$\int_{\Omega} k \nabla H(k - k_0) \cdot \nabla \phi \, d\Omega = k_0 \int_{\Omega} \nabla H(k - k_0) \cdot \nabla \phi \, d\Omega \quad \text{for all } \phi \in C_0^\infty.$$

Note that (4.14) ensures that $H(k - k_0) \neq 0$, and thus, if k satisfies (4.14) and (4.15) then $H(k(x, y) - k_0)$ is an eigenfunction, in the distributional sense, with associated eigenvalue $k_0 = k(x_0, y_0)$ for the generalized eigenvalue problem (4.1).

The next question is, of course, what happens if either assumption (4.14) or (4.15) fails to hold? This is the topic of Sections 4.2.2 and 4.2.3.

4.2.2 *Case II.* Let $k \in K$ and assume that $(x_0, y_0) \in \Omega$ is a point such that

$$k_x(x_0, y_0) = 0 \quad \text{and} \quad k_y(x_0, y_0) = 0. \quad (4.21)$$

In this case we will show that the Dirac delta function is an eigenfunction, in the distributional sense, for the generalized eigenvalue problem (4.12).

To this end, let δ denote the delta distribution associated with the point (x_0, y_0) , i.e., δ denotes the linear functional

$$\delta : C_0^\infty \rightarrow \mathbb{R}$$

such that the action of applying δ to $\phi \in C_0^\infty$ is given by

$$\langle \delta, \phi \rangle = \phi(x_0, y_0). \quad (4.22)$$

Note that we use the notation

$$\langle \delta, \phi \rangle = \int_{\Omega} \delta \phi \, d\Omega = \phi(x_0, y_0).$$

Recall the form (4.12) of the eigenvalue problem. Let $\phi \in C_0^\infty$ be an arbitrary test function and consider the integral

$$I_1 = \int_{\Omega} k \nabla \delta \cdot \nabla \phi \, d\Omega = \int_{\Omega} (k \delta_x \phi_x + k \delta_y \phi_y) \, d\Omega.$$

Now integration by parts implies that

$$\begin{aligned} I_1 &= - \int_{\Omega} (\delta [k_x \phi_x + k \phi_{xx}] + \delta [k_y \phi_y + k \phi_{yy}]) \, d\Omega \\ &= - [k_x(x_0, y_0) \phi_x(x_0, y_0) + k(x_0, y_0) \phi_{xx}(x_0, y_0) + k_y(x_0, y_0) \phi_y(x_0, y_0) + k(x_0, y_0) \phi_{yy}(x_0, y_0)] \\ &= - [k(x_0, y_0) \phi_{xx}(x_0, y_0) + k(x_0, y_0) \phi_{yy}(x_0, y_0)], \end{aligned}$$

where the last equality follows from assumption (4.21). On the other hand, by inserting $\lambda = k_0 = k(x_0, y_0)$ and $u = \delta$ in the right-hand side of (4.12) we find that

$$\begin{aligned} I_2 &= k_0 \int_{\Omega} \nabla \delta \cdot \nabla \phi \, d\Omega = k_0 \int_{\Omega} (\delta_x \phi_x + \delta_y \phi_y) \, d\Omega \\ &= -k_0 \int_{\Omega} (\delta \phi_{xx} + \delta \phi_{yy}) \, d\Omega \\ &= -k(x_0, y_0) [\phi_{xx}(x_0, y_0) + \phi_{yy}(x_0, y_0)]. \end{aligned}$$

Thus $I_1 = I_2$ and it follows that the δ -distribution associated with the point (x_0, y_0) is an ‘eigenfunction’ with eigenvalue $k(x_0, y_0)$, provided that (4.21) holds.

4.2.3 *Case III.* Let $k \in K$ and assume that $(x_0, y_0) \in \Omega$ is a point such that

$$|S(k, (x_0, y_0))| > 0 \quad (4.23)$$

(cf. (4.10) and (4.11)). This is, in fact, the simplest case. According to the definition of K in (4.11), $S(k, (x_0, y_0))$ contains an open and connected subset G with strictly positive measure, i.e., $|G| > 0$. This

ensures the existence of nonzero functions whose support is contained in G . Let u be such a function, i.e., we assume that

$$\text{supp}(u) \subset G. \quad (4.24)$$

Since $k(x, y) = k(x_0, y_0)$ on G , it follows that

$$\begin{aligned} \int_{\Omega} k \nabla u \cdot \nabla \phi \, d\Omega &= k(x_0, y_0) \int_G \nabla u \cdot \nabla \phi \, d\Omega \\ &= k(x_0, y_0) \int_{\Omega} \nabla u \cdot \nabla \phi \, d\Omega, \end{aligned}$$

which should be compared with (4.12). Hence we conclude that in this case every nonzero function satisfying (4.24) is a generalized eigenfunction with associated eigenvalue $k(x_0, y_0)$.

THEOREM 4.1 Let k be a coefficient function in the set K defined in (4.11). For every $(x_0, y_0) \in \Omega$ there exists a generalized function u such that $\lambda = k(x_0, y_0)$ and u forms an eigenvalue–eigenfunction pair for the generalized eigenvalue problem (4.1).

Furthermore, the following statements hold.

- If conditions (4.14) and (4.15) hold then

$$u = H(k - k_0)$$

is an eigenfunction with associated eigenvalue $\lambda = k_0 = k(x_0, y_0)$. Here H denotes the Heaviside function (see (4.16)).

- If condition (4.21) holds then

$$u = \delta_{(x_0, y_0)}$$

is a generalized eigenfunction with associated eigenvalue $\lambda = k_0 = k(x_0, y_0)$. Here $\delta_{(x_0, y_0)}$ is the Dirac delta distribution associated with the point (x_0, y_0) (see (4.22)).

- If condition (4.23) holds then any function u satisfying (4.24), where G is the set defined in (4.11), is a solution of the generalized eigenvalue problem (4.1) with associated eigenvalue $\lambda = k_0 = k(x_0, y_0)$.

5. Conclusions

In this paper we have analysed the eigenvalues and eigenfunctions of second-order elliptic PDEs preconditioned by the inverse of the Laplacian. We have shown by numerical experiments and mathematical analysis that there is a strong relation between the spectrum of the preconditioned operator and the range of the coefficient function k , provided that k is smooth and satisfies certain measure-theoretical properties.

More precisely, in the discrete case the spectrum seems to be accurately approximated by the values of the coefficient function evaluated at the mesh points. Furthermore, we have proven, both for the associated operators defined on Sobolev spaces and in terms of generalized functions, that the range of k is contained in the spectrum of the preconditioned operator.

The purpose of this paper has been to obtain a deeper understanding of the convergence properties of the CG method applied to second-order elliptic equations. For problems with large jump discontinuities in the coefficients the success of efficient preconditioners is commonly based on some sort of clustering effect of the eigenvalues. The present work shows that such an approach might be very difficult to apply to problems involving continuous coefficients with large variations. In particular, if the inverse Laplacian is applied as a preconditioner then such a clustering effect will not occur.

Acknowledgements

We would like to thank Kent-Andre Mardal for helping us with the implementation of the C++ software used in this work. Furthermore, we are very grateful to the referees for their most interesting comments and suggestions.

REFERENCES

- AXELSSON, O. (1994) *Iterative Solution Methods*. Cambridge: Cambridge University Press.
- AXELSSON, O. & LINDSKOG, G. (1986a) On the eigenvalue distribution of a class of preconditioning methods. *Numer. Math.*, **48**, 479–498.
- AXELSSON, O. & LINDSKOG, G. (1986b) On the rate of convergence of the preconditioned conjugate gradient method. *Numer. Math.*, **48**, 499–523.
- BEBENDORF, M. & HACKBUSCH, W. (2003) Existence of H -matrix approximants to the inverse FE matrix of elliptic operators with L^∞ -coefficients. *Numer. Math.*, **95**, 1–28.
- BECKERMANN, B. & KUIJLAARS, A. B. J. (2001) Superlinear convergence of conjugate gradients. *SIAM J. Numer. Anal.*, **39**, 300–329.
- BECKERMANN, B. & KUIJLAARS, A. B. J. (2002) Superlinear CG convergence for special right-hand sides. *Electron. Trans. Numer. Anal.*, **14**, 1–19. Special volume on orthogonal polynomials, approximation theory, and harmonic analysis (Inzell, 2000).
- BRAMBLE, J. H., PASCIAK, J. E., WANG, J. & XU, J. (1991) Convergence estimates for multigrid algorithms without regularity assumptions. *Math. Comp.*, **57**, 23–45.
- CAI, X., NIELSEN, B. F. & TVEITO, A. (1999) An analysis of a preconditioner for the discretized pressure equation arising in reservoir simulation. *IMA J. Numer. Anal.*, **19**, 291–316.
- CHAN, T. F. & MATHEW, T. P. (1994) Domain decomposition algorithms. *Acta Numer.*, 61–143.
- DRYJA, M. (1994) Multilevel methods for elliptic problems with discontinuous coefficients in three dimensions. *Domain Decomposition Methods in Scientific and Engineering Computing* (D. E. Keyes & J. Xu eds). Providence, RI: American Mathematical Society, pp. 43–47.
- DRYJA, M., SARKIS, M. & WIDLUND, O. B. (1996) Multilevel Schwarz methods for elliptic problems with discontinuous coefficients in three dimensions. *Numer. Math.*, **72**, 313–348.
- GRAHAM, I. G. & HAGGER, M. J. (1999) Unstructured additive Schwarz-conjugate gradient method for elliptic problems with highly discontinuous coefficients. *SIAM J. Sci. Comput.*, **20**, 2041–2066.
- GRIFFEL, D. H. (1981) *Applied Functional Analysis*. Chichester: Horwood.
- HACKBUSCH, W. (1994) *Iterative Solution of Large Sparse Systems of Equations*. New York: Springer.
- LANGTANGEN, H. P. (1999) *Computational Partial Differential Equations. Numerical Methods and Diffpack Programming*. Berlin: Springer.
- MARTI, J. T. (1986) *Introduction to Sobolev Spaces and Finite Element Solution of Elliptic Boundary Value Problems*. London: Academic Press.
- NAIMAN, A. E., BABUŠKA, I. M. & ELMAN, H. C. (1997) A note on conjugate gradient convergence. *Numer. Math.*, **76**, 209–230.
- NAIMAN, A. E. & ENGELBERG, S. (2000) A note on conjugate gradient convergence. II, III. *Numer. Math.*, **85**, 665–683, 685–696.
- ROYDEN, H. L. (1989) *Real Analysis*. New York: Macmillan.
- STOER, J. & BULIRSCH, R. (1993) *Introduction to Numerical Analysis*. New York: Springer.

Appendix A. A measure-theoretical lemma

Let $k : \Omega \rightarrow \mathbb{R}$ be a measurable function, let (x_0, y_0) be an arbitrary point in Ω and consider the sets

$$S_0 = \{(x, y) \in \Omega; k(x, y) = k(x_0, y_0)\}, \quad (\text{A.1})$$

$$S_\varepsilon = \{(x, y) \in \Omega; |k(x, y) - k(x_0, y_0)| \leq \varepsilon\}. \quad (\text{A.2})$$

Note that

$$S_0 \subset S_\varepsilon \quad \text{for all } \varepsilon > 0, \quad (\text{A.3})$$

and, furthermore,

$$\bigcap_{\varepsilon > 0} S_\varepsilon = S_0.$$

Motivated by these properties, we now establish the following lemma.

LEMMA A.1 Let $k : \Omega \rightarrow \mathbb{R}$ be a measurable function and let S_0 and S_ε denote the sets defined in (A.1) and (A.2), respectively. If the area of Ω is finite, i.e., $|\Omega| < \infty$, then

$$\lim_{\varepsilon \rightarrow 0} |S_\varepsilon| = |S_0|,$$

where $|S_0|$ and $|S_\varepsilon|$ denote the measures of these sets.

Proof. Let \mathcal{X}_0 and \mathcal{X}_ε denote the characteristic functions of S_0 and S_ε , respectively, that is,

$$\mathcal{X}_0(x, y) = \begin{cases} 1 & \text{for } (x, y) \in S_0, \\ 0 & \text{elsewhere} \end{cases}$$

and

$$\mathcal{X}_\varepsilon(x, y) = \begin{cases} 1 & \text{for } (x, y) \in S_\varepsilon, \\ 0 & \text{elsewhere.} \end{cases}$$

We start by proving that \mathcal{X}_ε converges point-wise toward \mathcal{X}_0 as $\varepsilon \rightarrow 0$. Let (\tilde{x}, \tilde{y}) be an arbitrary point in Ω .

(a) If $(\tilde{x}, \tilde{y}) \in S_0$ then (A.3) implies that

$$\mathcal{X}_0(\tilde{x}, \tilde{y}) = 1 = \mathcal{X}_\varepsilon(\tilde{x}, \tilde{y}) \quad \text{for all } \varepsilon > 0.$$

(b) If $(\tilde{x}, \tilde{y}) \notin S_0$ then

$$\mathcal{X}_0(\tilde{x}, \tilde{y}) = 0.$$

Moreover, for $\varepsilon < |k(\tilde{x}, \tilde{y}) - k(x_0, y_0)|$ it follows from the definition (A.2) of S_ε that $(\tilde{x}, \tilde{y}) \notin S_\varepsilon$. Hence

$$\mathcal{X}_\varepsilon(\tilde{x}, \tilde{y}) = 0 \quad \text{for all } \varepsilon < |k(\tilde{x}, \tilde{y}) - k(x_0, y_0)|.$$

Clearly, (a) and (b) imply that

$$\lim_{\varepsilon \rightarrow 0} \mathcal{X}_\varepsilon(\tilde{x}, \tilde{y}) = \mathcal{X}_0(\tilde{x}, \tilde{y}),$$

and we conclude that \mathcal{X}_ε converges point-wise toward \mathcal{X}_0 as $\varepsilon \rightarrow 0$.

Note that (cf. (A.3))

$$|1 - \mathcal{X}_\varepsilon| = 1 - \mathcal{X}_\varepsilon \leq 1 - \mathcal{X}_0 \quad \text{for all } \varepsilon > 0.$$

Thus the dominated convergence theorem (see, e.g., Royden, 1989) implies that

$$\lim_{\varepsilon \rightarrow 0} \int_{\Omega} (1 - \mathcal{X}_\varepsilon) \, d\Omega = \int_{\Omega} (1 - \mathcal{X}_0) \, d\Omega.$$

Consequently, since Ω has a finite measure,

$$\lim_{\varepsilon \rightarrow 0} |S_\varepsilon| = |S_0|,$$

which finishes the proof. □