

NUMERICAL ANALYSIS AND SCIENTIFIC COMPUTING
PREPRINT SERIA

**A robust preconditioner for high-contrast
problems**

Y. GORB

D. KURZANOVA

Y. KUZNETSOV

PREPRINT #58



DEPARTMENT OF MATHEMATICS
UNIVERSITY OF HOUSTON

JANUARY 2017

A Robust Preconditioner for High-Contrast Problems

Yuliya Gorb^{*1}, Daria Kurzanova^{†1}, and Yuri Kuznetsov^{‡1}

¹Department of Mathematics, University of Houston, Houston, TX 77204

Abstract

This paper concerns robust numerical treatment of an elliptic PDE with high contrast coefficients. A finite-element discretization of such an equation yields a linear system whose conditioning worsens as the variations in the values of PDE coefficients becomes large. This paper introduces a description of the problem whose discretization results in a linear system of saddle point type. Then a robust preconditioner for the Lancsoz method of minimized iterations used for solving the derived saddle point problem is proposed. Numerical examples demonstrate effectiveness and robustness of the proposed class of preconditioners yielding the number of iterations independent of the contrast and the discretization size.

Keywords: high contrast, saddle point problem, robust preconditioning, Schur complement, Lancsoz method

1 Introduction

In this paper, we consider the iterative solution of the linear system arising from the discretization of the diffusion problem

$$-\nabla \cdot [\sigma(x)\nabla u] = f, \quad x \in \Omega \quad (1)$$

with appropriate boundary conditions on $\Gamma = \partial\Omega$. We assume that Ω is a bounded domain $\Omega \subset \mathbb{R}^d$, $d \in \{2, 3\}$, that contains $m \geq 1$ polygonal or polyhedral subdomains \mathcal{D}^i , see Fig. 1. The main focus of this work is on the case when the coefficient function $\sigma(x) \in L^\infty(\Omega)$ varies largely within the domain Ω , that is,

$$\kappa = \frac{\sup_{x \in \Omega} \sigma(x)}{\inf_{x \in \Omega} \sigma(x)} \gg 1.$$

In this work, we assume that the domain Ω contains disjoint polygonal or polyhedral subdomains \mathcal{D}^i , $i \in \{1, \dots, m\}$, where σ takes “large” values, e.g. of order $O(\kappa)$, but remains of $O(1)$ in the domain outside of $\mathcal{D} := \cup_{i=1}^m \mathcal{D}^i$.

The P1-FEM discretization of this problem results in a linear system

$$\mathcal{A}\mathbf{x} = \mathbf{f}, \quad (2)$$

with a large and sparse matrix \mathcal{A} . A major issue in numerical treatments of (1) with the discussed above coefficient σ is that high contrast leads to an ill-conditioned matrix \mathcal{A} in (2). Indeed, if

*gorb@math.uh.edu

†dariak@math.uh.edu

‡yuri@math.uh.edu

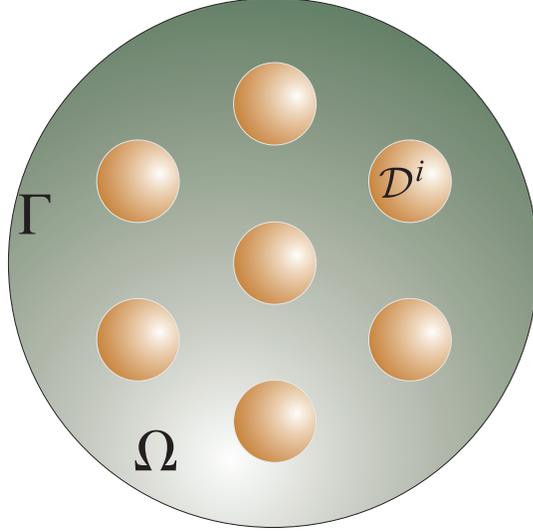


Figure 1: The domain Ω with highly conducting inclusions \mathcal{D}^i , $i \in \{1, \dots, m\}$

h is the discretization scale, then the condition number of the resulting stiffness matrix \mathcal{A} grows proportionally to h^{-2} with coefficient of proportionality depending on κ . Because of that, the high contrast problems have been a subject of an active research recently, see e.g. [1, 2].

There is one more feature of the system (2) that we investigate in this paper. Recall, that if \mathcal{A} is symmetric and positive definite, then (2) is typically solved with the Conjugate Gradient (CG) method, if \mathcal{A} is nonsymmetric then the most common solver for (2) is GMRES. Here, we focus on the type of continuum problems whose discrete approximation (2) yields a symmetric but indefinite matrix \mathcal{A} . In particular, we obtain a *saddle point system* [5, 19], in which \mathcal{A} is written in the block form:

$$\mathcal{A} = \begin{bmatrix} \mathbf{A} & \mathbf{B}^T \\ \mathbf{B} & -\mathbf{C} \end{bmatrix}, \quad (3)$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric positive definite, $\mathbf{B} \in \mathbb{R}^{k \times n}$ is rank deficient, and $\mathbf{C} \in \mathbb{R}^{k \times k}$ is symmetric and positive semidefinite, so that the corresponding linear system is singular but consistent. Unfortunately, Krylov space iterative methods tend to converge very slowly when applied to systems with saddle point matrices and preconditioners are needed to achieve faster convergence.

The special case of (2) with (3) tackled in this paper is when $\mathbf{C} \equiv \mathbf{0}$. It has been extensively studied by many authors, when \mathcal{A} is *nonsingular*, in which case \mathbf{B} must be of full rank, see e.g. [11, 15] and references therein. The CG method that was mainly developed for the iterative solution of linear systems with symmetric definite matrices is not in general robust for systems with indefinite matrices, [21]. The *Lanczos algorithm* of minimized iterations does not have such a restriction and has been utilized in this paper. Below in the paper, we introduce a construction of a robust preconditioner for solving (2) by the Lanczos iterative scheme, that is, whose convergence rate is independent of the contrast parameter $\kappa \gg 1$ and the discretization size $h > 0$.

The rest of the paper is organized as follows. In Chapter 2 the mathematical problem formulation is presented and main results are stated. Chapter 3 discusses proofs of main results, and numerical results of the proposed procedure are given in Chapter 4. Conclusions are presented in Chapter 5. Proof an auxiliary fact is given in Appendix 6.

Acknowledgements. First two authors were supported by the NSF grant DMS-1350248.

2 Problem Formulation and Main Results

Consider an open, a bounded domain $\Omega \subset \mathbb{R}^d$, $d \in \{2, 3\}$ with piece-wise smooth boundary Γ , that contains $m \geq 1$ subdomains \mathcal{D}^i , see Fig.1. For simplicity, we assume that Ω and \mathcal{D}^i are polygons if $d = 2$ or polyhedra if $d = 3$. The union of \mathcal{D}^i is denoted by \mathcal{D} . In the domain Ω we consider the elliptic problem

$$\begin{cases} -\nabla \cdot [\sigma(x)\nabla u] = f, & x \in \Omega \\ u = 0, & x \in \Gamma \end{cases} \quad (4)$$

with the coefficient σ that largely varies inside the domain Ω . For simplicity of the presentation, we focus on the case when σ is a piecewise constant function given by

$$\sigma(x) = \begin{cases} 1, & x \in \Omega \setminus \overline{\mathcal{D}} \\ 1 + \frac{1}{\varepsilon_i}, & x \in \mathcal{D}_i, i \in \{1, \dots, m\} \end{cases} \quad (5)$$

with $\max_i \varepsilon_i \ll 1$. We also assume that the source term in (4) is $f \in L^2(\Omega)$.

2.1 Derivation of a Singular Saddle Point Problem

When performing a FEM discretization of (4) with (5), we choose the FEM space $V_h \subset H_0^1(\Omega)$ to be the space of linear finite-element functions defined on a conforming quasi-uniform triangulation Ω_h of Ω of the size $h \ll 1$. For simplicity, we assume that $\partial\Omega_h = \Gamma$. If $\mathcal{D}_h^i = \Omega_h|_{\mathcal{D}^i}$ then we denote $V_h^i := V_h|_{\mathcal{D}_h^i}$ and $\mathcal{D}_h := \cup_{i=1}^m \mathcal{D}_h^i$. Then the FEM formulation of the problem (4)-(5) is to find $u_h \in V_h$ and $\lambda_h = (\lambda_h^1, \dots, \lambda_h^m)$ with $\lambda_h^i \in V_h^i$ such that

$$\int_{\Omega_h} \nabla u_h \cdot \nabla v_h \, dx + \int_{\mathcal{D}_h} \nabla \lambda_h \cdot \nabla v_h \, dx = \int_{\Omega_h} f v_h \, dx, \quad \forall v_h \in V_h, \quad (6)$$

provided

$$u_h = \varepsilon_i \lambda_h^i + c_i \quad \text{in } \mathcal{D}_h^i, \quad i \in \{1, \dots, m\}, \quad (7)$$

where c_i is an arbitrary constant. The FEM discretization of the problem (6) yields a system of linear equations

$$\mathbf{A}\bar{u} + \mathbf{B}^T \bar{\lambda} = \bar{F}. \quad (8)$$

Before providing the description of all elements of (8), we first introduce the following notations for the number of nodes in different parts of Ω_h . Let N be the total number of nodes in Ω_h , and n be the number of nodes in $\overline{\mathcal{D}_h}$ so that

$$n = \sum_{i=1}^m n_i,$$

where n_i denotes the number of degrees of freedom in $\overline{\mathcal{D}_h^i}$, and, finally, n_0 is the number of nodes in $\Omega_h \setminus \overline{\mathcal{D}_h}$, so that we have

$$N = n_0 + n = n_0 + \sum_{i=1}^m n_i.$$

Then in (8), the vector $\bar{u} \in \mathbb{R}^N$ has entries $u_i = u_h(x_i)$ with $x_i \in \bar{\Omega}_h$. We count the entries of \bar{u} in such a way that its first n elements correspond to the nodes of $\bar{\mathcal{D}}_h$, and the remaining n_0 entries correspond to the nodes of $\bar{\Omega}_h \setminus \bar{\mathcal{D}}_h$. Similarly, the vector $\bar{\lambda} \in \mathbb{R}^n$ has entries $\lambda_i = \lambda_h(x_i)$ where $x_i \in \bar{\mathcal{D}}_h$.

The symmetric positive definite matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ of (8) is the stiffness matrix that arise from the discretization of the Laplace operator with the homogeneous Dirichlet boundary conditions on Γ . Entries of \mathbf{A} are defined by

$$(\mathbf{A}\bar{u}, \bar{v}) = \int_{\Omega_h} \nabla u_h \cdot \nabla v_h \, dx, \quad \text{where } \bar{u}, \bar{v} \in \mathbb{R}^N, \quad u_h, v_h \in V_h, \quad (9)$$

where (\cdot, \cdot) is the standard dot-product of vectors. This matrix can also be partitioned into

$$\mathbf{A} = \begin{bmatrix} A_{\mathcal{D}\mathcal{D}} & A_{\mathcal{D}0} \\ A_{0\mathcal{D}} & A_{00} \end{bmatrix}, \quad (10)$$

where the block $A_{\mathcal{D}\mathcal{D}} \in \mathbb{R}^{n \times n}$ is the stiffness matrix corresponding to the highly conducting inclusions $\bar{\mathcal{D}}_h^i$, $i \in \{1, \dots, m\}$, the block $A_{00} \in \mathbb{R}^{n_0 \times n_0}$ corresponds to the region outside of $\bar{\mathcal{D}}_h$, and the entries of $A_{\mathcal{D}0} \in \mathbb{R}^{n \times n_0}$ and $A_{0\mathcal{D}} = A_{\mathcal{D}0}^T$ are assembled from contributions both from finite elements in $\bar{\mathcal{D}}_h$ and $\bar{\Omega}_h \setminus \bar{\mathcal{D}}_h$.

The matrix $\mathbf{B} \in \mathbb{R}^{n \times n}$ of (8) is also written in the block form as

$$\mathbf{B} = [\mathbf{B}_{\mathcal{D}} \quad \mathbf{0}] \quad (11)$$

with zero-matrix $\mathbf{0} \in \mathbb{R}^{n \times n_0}$ and $\mathbf{B}_{\mathcal{D}} \in \mathbb{R}^{n \times n}$ that corresponds to the highly conducting inclusions. The matrix $\mathbf{B}_{\mathcal{D}}$ is the stiffness matrix corresponding to the discretization of the Laplace operator in the domain $\bar{\mathcal{D}}_h$ with the Neumann boundary conditions on $\partial\mathcal{D}_h$. In its turn, $\mathbf{B}_{\mathcal{D}}$ is written in the block form by

$$\mathbf{B}_{\mathcal{D}} = \begin{bmatrix} \mathbf{B}_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \mathbf{B}_m \end{bmatrix} = \text{diag} (\mathbf{B}_1, \dots, \mathbf{B}_m)$$

with matrices $\mathbf{B}_i \in \mathbb{R}^{n_i \times n_i}$, whose entries are similarly defined by

$$(\mathbf{B}_i \bar{u}, \bar{v}) = \int_{\mathcal{D}_h^i} \nabla u_h \cdot \nabla v_h \, dx, \quad \text{where } \bar{u}, \bar{v} \in \mathbb{R}^{n_i}, \quad u_h, v_h \in V_h^i. \quad (12)$$

We remark that each \mathbf{B}_i is positive semidefinite with

$$\ker \mathbf{B}_i = \text{span} \left\{ \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \right\}. \quad (13)$$

Finally, the vector $\bar{\mathbf{F}} \in \mathbb{R}^n$ of (8) is defined in a similar way by

$$(\bar{\mathbf{F}}, \bar{v}) = \int_{\Omega_h} f v_h \, dx, \quad \text{where } \bar{v} \in \mathbb{R}^n, \quad v_h \in V_h.$$

To complete the derivation of the linear system corresponding to (6)-(7) we add the discrete analog of the relation (7). For that, denote

$$\Sigma_\varepsilon = \begin{bmatrix} \varepsilon_1 \mathbf{B}_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \varepsilon_m \mathbf{B}_m \end{bmatrix} = \text{diag} (\varepsilon_1 \mathbf{B}_1, \dots, \varepsilon_m \mathbf{B}_m)$$

then (7) implies

$$\Sigma_\varepsilon \bar{\lambda} = \mathbf{B}\bar{u}, \quad (14)$$

that together with (8) yields

$$\begin{cases} \mathbf{A}\bar{u} + \mathbf{B}^T \bar{\lambda} = \bar{\mathbf{F}}, \\ \mathbf{B}\bar{u} - \Sigma_\varepsilon \bar{\lambda} = \bar{\mathbf{0}}, \end{cases} \quad \bar{u} \in \mathbb{R}^N, \quad \mathbb{R}^n \ni \bar{\lambda} \perp \ker \mathbf{B}_\mathcal{D}, \quad (15)$$

or

$$\mathcal{A}_\varepsilon \mathbf{x}_\varepsilon = \bar{\mathcal{F}}, \quad (16)$$

where

$$\mathcal{A}_\varepsilon = \begin{bmatrix} \mathbf{A} & \mathbf{B}^T \\ \mathbf{B} & -\Sigma_\varepsilon \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{\mathcal{D}\mathcal{D}} & \mathbf{A}_{\mathcal{D}0} & \mathbf{B}_\mathcal{D} \\ \mathbf{A}_{0\mathcal{D}} & \mathbf{A}_{00} & \mathbf{0} \\ \mathbf{B}_\mathcal{D} & \mathbf{0} & -\Sigma_\varepsilon \end{bmatrix}, \quad \mathbf{x}_\varepsilon = \begin{bmatrix} \bar{u} \\ \bar{\lambda} \end{bmatrix}, \quad \bar{\mathcal{F}} = \begin{bmatrix} \bar{\mathbf{F}} \\ \bar{\mathbf{0}} \end{bmatrix}. \quad (17)$$

This saddle point formulation (16)-(17) was first proposed in [13]. Clearly, there exists a unique solution $\bar{u} \in \mathbb{R}^N$ and $\mathbb{R}^n \ni \bar{\lambda} \perp \ker \mathbf{B}_\mathcal{D}$ of (16)-(17).

2.2 Discussions on the system (15)

A few remarks regarding the linear system (15) are in order.

Remark 1. *The linear system (16)-(17) is the **saddle-point problem** with symmetric and indefinite matrix \mathcal{A}_ε . In the traditional treatments of saddle-point problems, see [5], it is typically assumed that the matrix \mathbf{B} has full rank. In our case, the matrix \mathbf{B} has a **nonzero kernel** due to (13).*

The choice of our discrete problem formulation (16)-(17) with the *singular* matrix \mathcal{A}_ε is motivated by that fact we can fully characterize $\ker \mathbf{B}_\mathcal{D}$ and take advantage of the structure of the singular matrix $\mathbf{B}_\mathcal{D}$ in our construction below.

Remark 2. *One can also introduce another FEM formulation **equivalent** to (6)-(7) that results in a **nonsingular** linear system (16).*

To derive it, we replace (7) with

$$\int_{\mathcal{D}_h^i} \nabla u_h \cdot \nabla v_h^i \, dx - \sum_{i=1}^m \varepsilon_i \int_{\mathcal{D}_h^i} \nabla \lambda_h^i \cdot \nabla v_h^i \, dx = 0, \quad i \in \{1, \dots, m\} \quad \text{for all } v_h^i \in V_h^i. \quad (18)$$

If, in addition, we assume that

$$\int_{\mathcal{D}_h^i} \lambda_h^i \, dx = 0, \quad i \in \{1, \dots, m\},$$

to fix arbitrary constants in (7), then obtain

$$\int_{\mathcal{D}_h^i} \nabla u_h \cdot \nabla v_h^i \, dx - \sum_{i=1}^m \varepsilon_i \int_{\mathcal{D}_h^i} \nabla \lambda_h^i \cdot \nabla v_h^i \, dx - \alpha_i^2 \left[\int_{\mathcal{D}_h^i} \lambda_h^i \, dx \right] \left[\int_{\mathcal{D}_h^i} v_h^i \, dx \right] = 0, \quad \text{for all } v_h^i \in V_h^i, \quad (19)$$

with some coefficients α_i . Then (19) together with (6) yields a FEM formulation equivalent to (6)-(7) and results in a **nonsingular** saddle-point problem (16). This is because the newly added term in (19) with appropriately chosen coefficients α_i , $i \in \{1, \dots, m\}$, describes an orthogonal projector on the kernel of the functional introduced in (12).

Remark 3. Since (6)-(7) admits an equivalent FEM description (6), (19) that yields the non-singular linear system (16) and since $\det \mathcal{A}_\varepsilon$ and all minors of \mathcal{A}_ε are continuously differentiable functions of its entries, hence, of ε_i , $i \in \{1, \dots, m\}$, we can expand $\mathcal{A}_\varepsilon^{-1}$ into the Taylor series of $(\varepsilon_1, \dots, \varepsilon_m)$ in some neighborhood of $(0, \dots, 0)$. Keeping the zero-order term of this expansion only we have

$$\mathcal{A}_\varepsilon^{-1} = \mathcal{A}_0^{-1}(1 + o(1)), \quad \text{as } \varepsilon_i \ll 1, \quad i \in \{1, \dots, m\}. \quad (20)$$

Using the above asymptotic relation (20), below we construct a preconditioner for the matrix \mathcal{A}_0 , that we later use in the Lanczos algorithm when numerically solving (4).

Denote the solution of (16)-(17) by

$$\mathbf{x}_\varepsilon = \begin{bmatrix} \bar{u}_\varepsilon \\ \bar{\lambda}_\varepsilon \end{bmatrix},$$

and consider an auxiliary linear system

$$\mathcal{A}_0 \mathbf{x}_0 = \begin{bmatrix} \mathbf{A} & \mathbf{B}^T \\ \mathbf{B} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \bar{u}_0 \\ \bar{\lambda}_0 \end{bmatrix} = \begin{bmatrix} \bar{\mathbf{F}} \\ \bar{\mathbf{0}} \end{bmatrix}, \quad (21)$$

or

$$\begin{cases} \mathbf{A} \bar{u}_0 + \mathbf{B}^T \bar{\lambda}_0 &= \bar{\mathbf{F}}, \\ \mathbf{B} \bar{u}_0 &= \bar{\mathbf{0}}. \end{cases} \quad (22)$$

where matrices \mathbf{A} , \mathbf{B} and the vector $\bar{\mathbf{F}}$ are the same as above. The matrix \mathcal{A}_0 of (21) and the one of (20) is the same. The linear system (21) or, equivalently, (22) emerges in a FEM discretization of the diffusion problem posed in the domain Ω whose inclusions are *infinitely conducting*, that is, when $\varepsilon = 0$ in (5). The PDE formulation of this problem is as follows (see e.g. [7])

$$\begin{cases} \Delta u = f, & x \in \Omega \setminus \bar{\mathcal{D}} \\ u = \text{const}, & x \in \partial \mathcal{D}^i, \quad i \in \{1, \dots, m\} \\ \int_{\partial \mathcal{D}^i} \nabla u \cdot \mathbf{n}_i \, ds = 0, & i \in \{1, \dots, m\} \\ u = 0, & x \in \Gamma \end{cases} \quad (23)$$

where \mathbf{n}_i is the outer unit normal to the surface $\partial \mathcal{D}^i$. If $u \in H_0^1(\Omega \setminus \bar{\mathcal{D}})$ is an electric potential then it attains constant values on the inclusions \mathcal{D}^i and these constants are not known a priori so that they are unknowns of the problem (23) together with u .

Formulation (21) or (22) also arises in constrained quadratic optimization problem and solving the Stokes equations for an incompressible fluid [8], and solving elliptic problems using methods combining fictitious domain and distributed Lagrange multiplier techniques to force boundary conditions [9].

Then the following relation between solutions of systems (15) and (22) holds true.

Lemma 1. Let $\mathbf{x}_0 = \begin{bmatrix} \bar{u}_0 \\ \bar{\lambda}_0 \end{bmatrix} \in \mathbb{R}^{N+n}$ be the solution of the linear system (22), and $\mathbf{x}_\varepsilon = \begin{bmatrix} \bar{u}_\varepsilon \\ \bar{\lambda}_\varepsilon \end{bmatrix} \in \mathbb{R}^{N+n}$ the solution of (15). Then

$$\bar{u}_\varepsilon \rightarrow \bar{u}_0 \quad \text{as } \varepsilon \rightarrow 0.$$

This lemma asserts that the discrete approximation for the problem (4)-(5) converges to the discrete approximation of the solution of (23) as $\varepsilon \rightarrow 0$. We also note that the continuum version of this fact was shown in [7]. For the reader's convenience the proof of this lemma is posted in Appendices below.

2.3 Spectral Properties of the Matrix \mathcal{A}_0 of the Auxiliary Problem (22)

It was previously observed, see e.g. [12], that the following matrix

$$\mathbf{P} = \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{BA}^{-1}\mathbf{B}^T \end{bmatrix}, \quad (24)$$

is the best choice for a preconditioner of \mathcal{A}_0 . This is because there are exactly three eigenvalues of \mathcal{A}_0 associated with the following generalized eigenvalue problem

$$\mathcal{A}_0 \begin{bmatrix} \bar{u} \\ \bar{\lambda} \end{bmatrix} = \mu \mathbf{P} \begin{bmatrix} \bar{u} \\ \bar{\lambda} \end{bmatrix}, \quad (25)$$

and they are: $\mu_1 < 0$, $\mu_2 = 1$ and $\mu_3 > 1$, and, hence, a Krylov subspace iteration method applied for a preconditioned system for solving (25) with (24) *converges to the exact solution in three iterations*.

The preconditioner (24) is also the best choice for our original problem (16)-(17) with $\varepsilon > 0$ as the eigenvalue of the generalized eigenvalue problem

$$\mathcal{A}_\varepsilon \mathbf{x} = \mu \mathbf{P} \mathbf{x}$$

belong to the union of $[c_1, c_2] \cup [c_3, c_4]$ with $c_1 \leq c_2 < 0$ and $0 < c_3 \leq c_4$, with numbers c_i being dependent on eigenvalues of (25) but not h , see [13].

Due to expensive evaluation of \mathbf{A}^{-1} in (24) makes \mathbf{P} of limited practical use, so that \mathbf{P} is a subject of primarily theoretical interest. To construct a preconditioner that one can actually use in practice, we seek for a matrix

$$\mathcal{P} = \begin{bmatrix} \mathcal{P}_A & \mathbf{0} \\ \mathbf{0} & \mathcal{P}_B \end{bmatrix}, \quad (26)$$

such that there exist constants α, β independent on the mesh size h and that

$$\alpha(\mathbf{P}\mathbf{x}, \mathbf{x}) \leq (\mathcal{P}\mathbf{x}, \mathbf{x}) \leq \beta(\mathbf{P}\mathbf{x}, \mathbf{x}) \quad \text{for all } \mathbf{x} \in \mathbb{R}^N. \quad (27)$$

This property (27) is hereafter referred to as *spectral equivalence* of \mathcal{P} to \mathbf{P} of (24). Below, we will construct \mathcal{P} of the form (26) in such a way that the block \mathcal{P}_A is spectrally equivalent to \mathbf{A} , whereas \mathcal{P}_B to $\mathbf{BA}^{-1}\mathbf{B}^T$. For the former one we can use any existing preconditioner developed for symmetric and positive definite matrices. Our primary aim is to construct a preconditioner \mathcal{P}_B that could be effectively used in solving (15).

2.4 Main Result: Block-Diagonal Preconditioner

The main theoretical result of this paper establishes a robust preconditioner for solving (21) and is given in the following theorem.

Theorem 1. *Let the triangulation Ω_h for (23) be conforming and quasi-uniform. Then the matrix $\mathbf{B}_\mathcal{D}$ is spectrally equivalent to the matrix $\mathbf{BA}^{-1}\mathbf{B}^T$, that is, there exist constants $\mu_\star, \mu^\star > 0$ independent of h and such that*

$$\mu_\star \leq \frac{(\mathbf{B}_\mathcal{D} \bar{\psi}, \bar{\psi})}{(\mathbf{BA}^{-1}\mathbf{B}^T \bar{\psi}, \bar{\psi})} \leq \mu^\star, \quad \text{for all } 0 \neq \bar{\psi} \in \mathbb{R}^n, \bar{\psi} \perp \ker \mathbf{B}_\mathcal{D}. \quad (28)$$

This theorem asserts that the nonzero eigenvalues of the generalized eigenproblem

$$\mathbf{B}\mathbf{A}^{-1}\mathbf{B}^T\bar{\psi} = \mu\mathbf{B}_D\bar{\psi}, \quad \bar{\psi} \in \mathbb{R}^n, \quad (29)$$

are bounded. Hence, its proof is based on the construction of the **upper** and **lower** bounds for μ in (29) and is comprised of the following facts many of which are proven in the next section.

Lemma 2. *The following equality of matrices holds*

$$\mathbf{B}\mathbf{A}^{-1}\mathbf{B}^T = \mathbf{B}_D\mathbf{S}_{00}^{-1}\mathbf{B}_D^T, \quad (30)$$

where

$$\mathbf{S}_{00} = \mathbf{A}_{\mathcal{D}\mathcal{D}} - \mathbf{A}_{\mathcal{D}0}\mathbf{A}_{00}^{-1}\mathbf{A}_{0\mathcal{D}},$$

is the Schur complement to the block \mathbf{A}_{00} of the matrix \mathbf{A} of (21).

This fact is straightforward and comes from the block structure of matrices \mathbf{A} of (10) and \mathbf{B} of (11). Indeed, using this, the generalized eigenproblem (29) can be rewritten as

$$\mathbf{B}_D\mathbf{S}_{00}^{-1}\mathbf{B}_D\bar{\psi} = \mu\mathbf{B}_D\bar{\psi}, \quad \bar{\psi} \in \mathbb{R}^n. \quad (31)$$

Introduce a matrix $\mathbf{B}_D^{1/2}$ via $\mathbf{B}_D = \mathbf{B}_D^{1/2}\mathbf{B}_D^{1/2}$ and note that $\ker \mathbf{B}_D = \ker \mathbf{B}_D^{1/2}$.

Lemma 3. *The generalized eigenvalue problem (31) is equivalent to*

$$\mathbf{B}_D^{1/2}\mathbf{S}_{00}^{-1}\mathbf{B}_D^{1/2}\bar{\varphi} = \mu\bar{\varphi}, \quad (32)$$

in the sense that they both have the same eigenvalues μ 's, and the corresponding eigenvectors are related via $\bar{\varphi} = \mathbf{B}_D^{1/2}\bar{\psi} \in \mathbb{R}^n$.

Lemma 4. *The generalized eigenvalue problem (32) is equivalent to*

$$\mathbf{B}_D\bar{u}_D = \mu\mathbf{S}_{00}\bar{u}_D, \quad (33)$$

in the sense that both problems have the same eigenvalues μ 's, and the corresponding eigenvectors are related via $\bar{u}_D = \mathbf{S}_{00}^{-1}\mathbf{B}_D^{1/2}\bar{\varphi} \in \mathbb{R}^n$.

This result is also straightforward and can be obtained multiplying (32) by $\mathbf{S}_{00}^{-1}\mathbf{B}_D^{1/2}$.

To that end, establishing the upper and lower bounds for the eigenvalues of (33) and due to equivalence of (33) with (32), (31), we obtain that eigenvalues of (29) are bounded. We are interested in nonzero eigenvalues of (33) for which the following result holds.

Lemma 5. *Let the triangulation Ω_h for (23) be conforming and quasi-uniform. Then there exists $\hat{\mu}_\star > 0$ independent of the mesh size $h > 0$ such that*

$$\hat{\mu}_\star \leq \frac{(\mathbf{B}_D\bar{u}_D, \bar{u}_D)}{(\mathbf{S}_{00}\bar{u}_D, \bar{u}_D)} \leq 1, \quad \text{for all } 0 \neq \bar{u}_D \in \mathbb{R}^n, \quad \bar{u}_D \perp \ker \mathbf{B}_D. \quad (34)$$

3 Proofs of statements of Chapter 2.4

3.1 Harmonic extensions

We now recall some classical results from the theory of elliptic PDEs. Suppose a function $u^{\mathcal{D}} \in H^1(\mathcal{D})$, then consider its harmonic extension $u^0 \in H^1(\Omega \setminus \overline{\mathcal{D}})$ that satisfies

$$\begin{cases} -\Delta u^0 = 0, & \text{in } \Omega \setminus \overline{\mathcal{D}}, \\ u^0 = u^{\mathcal{D}}, & \text{on } \partial\mathcal{D}, \\ u^0 = 0, & \text{on } \Gamma. \end{cases} \quad (35)$$

For such functions the following holds true:

$$\int_{\Omega} |\nabla u|^2 \, dx = \min_{v \in H_0^1(\Omega)} \int_{\Omega} |\nabla v|^2 \, dx, \quad (36)$$

where

$$u = \begin{cases} u^{\mathcal{D}}, & \text{in } \mathcal{D} \\ u^0, & \text{in } \Omega \setminus \overline{\mathcal{D}} \end{cases} \quad \text{and} \quad v = \begin{cases} u^{\mathcal{D}}, & \text{in } \mathcal{D} \\ v^0, & \text{in } \Omega \setminus \overline{\mathcal{D}} \end{cases}$$

where the function $v^0 \in H^1(\Omega \setminus \overline{\mathcal{D}})$ such that $v^0|_{\Gamma} = 0$, and

$$\|u\|_{H_0^1(\Omega)} \leq C \|u^{\mathcal{D}}\|_{H^1(\mathcal{D})} \quad \text{with the constant } C \text{ independent of } u^{\mathcal{D}}. \quad (37)$$

In view of (36), the function u^0 of (36) is the *best extension* of $u^{\mathcal{D}} \in H^1(\mathcal{D})$ among all $H^1(\Omega \setminus \overline{\mathcal{D}})$ functions that vanish on Γ . The algebraic linear system that corresponds to (36) satisfies the similar property. Namely, if the vector $\bar{u}_0 \in \mathbb{R}^{n_0}$ is a FEM discretization of the function $u^0 \in H_0^1(\Omega \setminus \overline{\mathcal{D}})$ of (35), then for a given $\bar{u}_{\mathcal{D}} \in \mathbb{R}^n$, the best extension $\bar{u}_0 \in \mathbb{R}^{n_0}$ would satisfy

$$\mathbf{A}_{0\mathcal{D}} \bar{u}_{\mathcal{D}} + \mathbf{A}_{00} \bar{u}_0 = 0, \quad (38)$$

and

$$\left(\mathbf{A} \begin{bmatrix} \bar{u}_{\mathcal{D}} \\ \bar{u}_0 \end{bmatrix}, \begin{bmatrix} \bar{u}_{\mathcal{D}} \\ \bar{u}_0 \end{bmatrix} \right) = \min_{\bar{v}_0 \in \mathbb{R}^{n_0}} \left(\mathbf{A} \begin{bmatrix} \bar{u}_{\mathcal{D}} \\ \bar{v}_0 \end{bmatrix}, \begin{bmatrix} \bar{u}_{\mathcal{D}} \\ \bar{v}_0 \end{bmatrix} \right). \quad (39)$$

Hereafter, we will use the index \mathcal{D} to indicate vectors or functions associated with the domain \mathcal{D} that is the union of all inclusions, and index 0 to indicate quantities that are associated with the domain outside the inclusions $\Omega \setminus \overline{\mathcal{D}}$.

3.2 Proof of Lemma 3

Consider generalized eigenvalue problem (31) and replace $\mathbf{B}_{\mathcal{D}}$ with $\mathbf{B}_{\mathcal{D}}^{1/2} \mathbf{B}_{\mathcal{D}}^{1/2}$ there, then

$$\mathbf{B}_{\mathcal{D}}^{1/2} \mathbf{B}_{\mathcal{D}}^{1/2} \mathbf{S}_{00}^{-1} \mathbf{B}_{\mathcal{D}}^{1/2} \mathbf{B}_{\mathcal{D}}^{1/2} \bar{\psi} = \mu \mathbf{B}_{\mathcal{D}}^{1/2} \mathbf{B}_{\mathcal{D}}^{1/2} \bar{\psi}.$$

Now multiply both sides by the Moore-Penrose pseudo inverse¹ $[\mathbf{B}_{\mathcal{D}}^{1/2}]^{\dagger}$, see e.g. [3]:

$$[\mathbf{B}_{\mathcal{D}}^{1/2}]^{\dagger} \mathbf{B}_{\mathcal{D}}^{1/2} \mathbf{B}_{\mathcal{D}}^{1/2} \mathbf{S}_{00}^{-1} \mathbf{B}_{\mathcal{D}}^{1/2} \mathbf{B}_{\mathcal{D}}^{1/2} \bar{\psi} = \mu [\mathbf{B}_{\mathcal{D}}^{1/2}]^{\dagger} \mathbf{B}_{\mathcal{D}}^{1/2} \mathbf{B}_{\mathcal{D}}^{1/2} \bar{\psi}.$$

¹ \mathbf{M}^{\dagger} is the Moore-Penrose pseudo inverse of \mathbf{M} if and only if it satisfies the following Moore-Penrose equations:

(i) $\mathbf{M}^{\dagger} \mathbf{M} \mathbf{M}^{\dagger} = \mathbf{M}^{\dagger}$, (ii) $\mathbf{M} \mathbf{M}^{\dagger} \mathbf{M} = \mathbf{M}$, (iii) $\mathbf{M} \mathbf{M}^{\dagger}$ and $\mathbf{M}^{\dagger} \mathbf{M}$ are symmetric.

This pseudo inverse has the property that

$$\left[\mathbf{B}_{\mathcal{D}}^{1/2} \right]^\dagger \mathbf{B}_{\mathcal{D}}^{1/2} = \mathbf{P}_{\text{im}},$$

where \mathbf{P}_{im} is an orthogonal projector onto the image $\mathbf{B}_{\mathcal{D}}^{1/2}$, hence, $\mathbf{P}_{\text{im}} \mathbf{B}_{\mathcal{D}}^{1/2} = \mathbf{B}_{\mathcal{D}}^{1/2}$ and therefore,

$$\mathbf{B}_{\mathcal{D}}^{1/2} \mathbf{S}_{00}^{-1} \mathbf{B}_{\mathcal{D}}^{1/2} \bar{\varphi} = \mu \bar{\varphi}, \quad \text{where} \quad \bar{\varphi} = \mathbf{B}_{\mathcal{D}}^{1/2} \bar{\psi}.$$

Conversely, consider the eigenvalue problem (32), and multiply its both sides by $\mathbf{B}_{\mathcal{D}}^{1/2}$. Then

$$\mathbf{B}_{\mathcal{D}}^{1/2} \mathbf{B}_{\mathcal{D}}^{1/2} \mathbf{S}_{00}^{-1} \mathbf{B}_{\mathcal{D}}^{1/2} \bar{\varphi} = \mu \mathbf{B}_{\mathcal{D}}^{1/2} \bar{\varphi},$$

where we replace $\bar{\varphi}$ by $\mathbf{B}_{\mathcal{D}}^{1/2} \bar{\psi}$

$$\mathbf{B}_{\mathcal{D}}^{1/2} \mathbf{B}_{\mathcal{D}}^{1/2} \mathbf{S}_{00}^{-1} \mathbf{B}_{\mathcal{D}}^{1/2} \mathbf{B}_{\mathcal{D}}^{1/2} \bar{\psi} = \mu \mathbf{B}_{\mathcal{D}}^{1/2} \mathbf{B}_{\mathcal{D}}^{1/2} \bar{\psi}$$

to obtain (31). \square

3.3 Proof of Lemma 5

I. Upper Bound for the Generalized Eigenvalues of (29)

Consider $\bar{\mathbf{u}} = \begin{bmatrix} \bar{u}_{\mathcal{D}} \\ \bar{u}_0 \end{bmatrix} \in \mathbb{R}^N$ with $\bar{u}_{\mathcal{D}} \in \mathbb{R}^n$, $\bar{u}_{\mathcal{D}} \perp \ker \mathbf{B}_{\mathcal{D}}$, and $\bar{u}_0 \in \mathbb{R}^{n_0}$ satisfying (38), then

$$(\mathbf{S}_{00} \bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}}) = (\mathbf{A} \bar{\mathbf{u}}, \bar{\mathbf{u}}). \quad (40)$$

from which using (9) and (12) we obtain:

$$\mu = \frac{(\mathbf{B}_{\mathcal{D}} \bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}})}{(\mathbf{S}_{00} \bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}})} = \frac{(\mathbf{B}_{\mathcal{D}} \bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}})}{(\mathbf{A} \bar{\mathbf{u}}, \bar{\mathbf{u}})} = \frac{\int_{\mathcal{D}_h} |\nabla u_h^{\mathcal{D}}|^2 \, dx}{\int_{\Omega_h} |\nabla u_h|^2 \, dx} \leq 1, \quad (41)$$

with

$$u_h = \begin{cases} u_h^{\mathcal{D}}, & \text{in } \mathcal{D}_h \\ u_h^0, & \text{in } \Omega \setminus \bar{\mathcal{D}}_h \end{cases} \quad (42)$$

where u_h^0 is the harmonic extension of $u_h^{\mathcal{D}}$ into $\Omega_h \setminus \bar{\mathcal{D}}_h$ in the sense (35). \square

II. Lower Bound for the Generalized Eigenvalues of (29)

Before providing the proofs, we introduce one more construction to simplify our consideration below. Because all inclusions are located at distances that are comparable to their sizes, we construct new domains $\hat{\mathcal{D}}^i$, $i \in \{1, \dots, m\}$, see Fig. 2, centered at the centers of the original \mathcal{D}^i but of sizes much larger of those of \mathcal{D}^i and such that

$$\hat{\mathcal{D}}^i \cap \hat{\mathcal{D}}^j = \emptyset, \quad \text{for } i \neq j.$$

From it follows below, one can see that the problem (23) might be partitioned into m independent subproblems, with what, without loss of generality, we assume that there is only one inclusion, that is, $m = 1$.

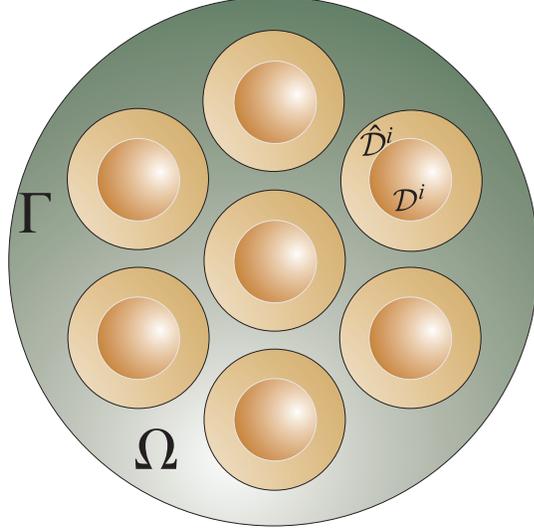


Figure 2: New domains $\hat{\mathcal{D}}^i$ for our construction of the lower bound of μ

We also recall a few important results from classical PDE theory analogs of which will be used below. Namely, for a given $v \in H^1(\mathcal{D})$ there exists an extension v_0 of v to $\Omega \setminus \overline{\mathcal{D}}$ so that

$$\|v_0\|_{H^1(\Omega \setminus \mathcal{D})} \leq C \|v\|_{H^1(\mathcal{D})}, \quad \text{with } C = C(d, \mathcal{D}, \Omega), \quad (43)$$

where $\|\cdot\|_{H^1(\Omega)}$ denotes the standard norm of $H^1(\Omega)$:

$$\|v\|_{H^1(\Omega)}^2 = \int_{\Omega} |\nabla v|^2 dx + \int_{\Omega} v^2 dx. \quad (44)$$

One can also introduce a number of norms equivalent to (44), and, in particular, below we will use

$$\|v\|_{\mathcal{D}}^2 := \int_{\mathcal{D}} |\nabla v|^2 dx + \frac{1}{R^2} \int_{\mathcal{D}} v^2 dx, \quad (45)$$

where R is the radius of the particle $\mathcal{D} = \mathcal{D}_1$. The scaling factor $1/R^2$ is needed for transforming the classical results from a reference (i.e. unit) disk to the disk of radius $R \neq 1$.

We note that the FEM analog of the extension result of (43) for a regular grid was shown in [20], from which it also follows that the constant C of (43) is independent of the mesh size h . We utilize this observation in our construction below.

Consider $u_h \in V_h$ given by (42). Introduce a space $\hat{V}_h = \{v_h \in V_h : v_h = 0 \text{ in } \Omega_h \setminus \overline{\mathcal{D}}_h\}$. Similarly to (42), define

$$\hat{V}_h \ni \hat{u}_h = \begin{cases} u_h^{\mathcal{D}}, & \text{in } \mathcal{D}_h \\ \hat{u}_h^0, & \text{in } \Omega_h \setminus \overline{\mathcal{D}}_h \end{cases}, \quad (46)$$

where \hat{u}_h^0 is the harmonic extension of $u_h^{\mathcal{D}}$ into $\hat{\mathcal{D}}_h \setminus \overline{\mathcal{D}}_h$ in the sense (35) and $\hat{u}_h^0 = 0$ on $\partial \hat{\mathcal{D}}_h$. Also, by (36) we have

$$\int_{\Omega_h \setminus \mathcal{D}_h} |\nabla u_h^0|^2 dx \leq \int_{\Omega_h \setminus \mathcal{D}_h} |\nabla \hat{u}_h^0|^2 dx.$$

Define the matrix

$$\hat{\mathbf{A}} := \begin{bmatrix} A_{\mathcal{D}\mathcal{D}} & \hat{A}_{\mathcal{D}0} \\ \hat{A}_{0\mathcal{D}} & \hat{A}_{00} \end{bmatrix}$$

by

$$\left(\hat{\mathbf{A}} \bar{v}, \bar{w} \right) = \int_{\Omega_h} \nabla v_h \cdot \nabla w_h dx, \quad \text{where } \bar{v}, \bar{w} \in \mathbb{R}^N, \quad v_h, w_h \in \hat{V}_h.$$

As before, introduce the Schur complement to the block \hat{A}_{00} of $\hat{\mathbf{A}}$:

$$\hat{S}_{00} = A_{\mathcal{D}\mathcal{D}} - \hat{A}_{\mathcal{D}0} \hat{A}_{00}^{-1} \hat{A}_{0\mathcal{D}}, \quad (47)$$

and consider a new generalized eigenvalue problem

$$B_{\mathcal{D}} \bar{u}_{\mathcal{D}} = \hat{\mu} \hat{S}_{00} \bar{u}_{\mathcal{D}} \quad \text{with} \quad \mathbb{R}^n \ni \bar{u}_{\mathcal{D}} \perp \ker B_{\mathcal{D}}. \quad (48)$$

By (39) and (40) we have

$$(S_{00} \bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}}) \leq \left(\hat{S}_{00} \bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}} \right) \quad \text{for all} \quad \bar{u}_{\mathcal{D}} \in \mathbb{R}^n. \quad (49)$$

Now, we consider a new generalized eigenvalue problem similar to one in (32), namely,

$$B_{\mathcal{D}}^{1/2} \hat{S}_{00}^{-1} B_{\mathcal{D}}^{1/2} \bar{\varphi} = \hat{\mu} \bar{\varphi}, \quad \bar{\varphi} \in \mathbb{R}^n. \quad (50)$$

We plan to replace $B_{\mathcal{D}}^{1/2}$ in (50) with a new symmetric positive-definite matrix $\hat{B}_{\mathcal{D}}^{1/2}$ so that

$$B_{\mathcal{D}}^{1/2} B_{\mathcal{D}}^{1/2} \bar{\xi} = B_{\mathcal{D}}^{1/2} \hat{B}_{\mathcal{D}}^{1/2} \bar{\xi} = \hat{B}_{\mathcal{D}}^{1/2} B_{\mathcal{D}}^{1/2} \bar{\xi} \quad \text{for all} \quad \mathbb{R}^n \ni \bar{\xi} \perp \ker B_{\mathcal{D}}, \quad (51)$$

with what (50) has the same nonzero eigenvalues as the problem

$$\hat{B}_{\mathcal{D}}^{1/2} \hat{S}_{00}^{-1} \hat{B}_{\mathcal{D}}^{1/2} \bar{\varphi} = \hat{\mu} \bar{\varphi}, \quad \bar{\varphi} \in \mathbb{R}^n. \quad (52)$$

For this purpose, we consider the decomposition:

$$B_{\mathcal{D}} = W \Lambda W^T,$$

where $W \in \mathbb{R}^{n \times n}$ is an orthogonal matrix composed of eigenvectors \bar{w}_i , $i \in \{0, 1, \dots, n-1\}$, of

$$B_{\mathcal{D}} \bar{w} = \nu \bar{w}, \quad \bar{w} \in \mathbb{R}^n,$$

and

$$\Lambda = \text{diag} [\nu_0, \nu_1, \dots, \nu_{n-1}].$$

Then \bar{w}_0 is an eigenvector of $B_{\mathcal{D}}$ corresponding to $\nu_0 = 0$ and

$$\bar{w}_0 = \frac{1}{\sqrt{n}} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}.$$

To that end, we choose

$$\hat{B}_{\mathcal{D}} = B_{\mathcal{D}} + \beta \bar{w}_0 \otimes \bar{w}_0 = B_{\mathcal{D}} + \beta \bar{w}_0 \bar{w}_0^T, \quad (53)$$

where $\beta > 0$ is some constant parameter chosen below. Note that the matrix $\hat{\mathbf{B}}_{\mathcal{D}}$ is symmetric and positive-definite, and satisfies (51). It is trivial to show that $\hat{\mathbf{B}}_{\mathcal{D}}$ given by (53) is spectrally equivalent to $\mathbf{B}_{\mathcal{D}} + \beta\mathbf{I}$ for any $\beta > 0$. Also, for quasi-uniform grids, the matrix $h^2\mathbf{I}$ (in 3-dim case, $h^3\mathbf{I}$) is spectrally equivalent to the mass matrix $\mathbf{M}_{\mathcal{D}}$, see e.g. [17], that implies there exists a constant $C > 0$ independent of h such that

$$\left(\hat{\mathbf{B}}_{\mathcal{D}}\bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}}\right) \geq C \left(\left(\mathbf{B}_{\mathcal{D}} + \frac{1}{R^2}\mathbf{M}_{\mathcal{D}} \right) \bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}} \right), \quad \text{with } \beta = \frac{h^2}{R^2}. \quad (54)$$

The choice of the matrix $\mathbf{B}_{\mathcal{D}} + \frac{1}{R^2}\mathbf{M}_{\mathcal{D}}$ for the spectral equivalence was motivated by the fact that the right hand side of (54) describes $\|\cdot\|_{\mathcal{D}_h}$ -norm (45) of the FEM function $u_h^{\mathcal{D}} \in V_h^1$ that corresponds to the vector $\bar{u}_{\mathcal{D}} \in \mathbb{R}^n$.

Now consider $\bar{u} = \begin{bmatrix} \bar{u}_{\mathcal{D}} \\ \bar{u}_0 \end{bmatrix} \in \mathbb{R}^N$ with $\bar{u}_{\mathcal{D}} \in \mathbb{R}^n$, $\bar{u}_{\mathcal{D}} \perp \ker \mathbf{B}_{\mathcal{D}}$, and $\bar{u}_0 \in \mathbb{R}^{n_0}$ satisfying (38), and similarly choose $\hat{u} = \begin{bmatrix} \bar{u}_{\mathcal{D}} \\ \hat{u}_0 \end{bmatrix} \in \mathbb{R}^N$ with $\hat{u}_0 \in \mathbb{R}^{n_0}$ satisfying $\hat{\mathbf{A}}_{0\mathcal{D}} \bar{u}_{\mathcal{D}} + \hat{\mathbf{A}}_{00} \bar{u}_0 = 0$, which implies

$$\left(\hat{\mathbf{S}}_{00}\bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}}\right) = \left(\hat{\mathbf{A}}\hat{u}, \hat{u}\right). \quad (55)$$

Then

$$\left(\hat{\mathbf{A}}\hat{u}, \hat{u}\right) = \int_{\Omega_h} |\nabla \hat{u}_h|^2 dx = \int_{\hat{\mathcal{D}}_h \setminus \mathcal{D}_h} |\nabla \hat{u}_h^0|^2 dx + \int_{\mathcal{D}_h} |\nabla u_h^{\mathcal{D}}|^2 dx \leq (C^* + 1) \|u_h^{\mathcal{D}}\|_{\mathcal{D}_h}^2, \quad (56)$$

where $\hat{u}_h \in \hat{V}_h$ is the same extension of $u_h^{\mathcal{D}}$ from $\bar{\mathcal{D}}_h$ to $\Omega_h \setminus \bar{\mathcal{D}}_h$ as defined in (46). For the inequality of (56), we applied the FEM analog of the extension result of (43) by [20], that yields that the constant C^* in (56) is independent of h .

With all the above, we have the following chain of inequalities:

$$\begin{aligned} & \frac{(\mathbf{B}_{\mathcal{D}}\bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}})}{(\mathbf{S}_{00}\bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}})} \stackrel{(51),(53)}{=} \frac{((\mathbf{B}_{\mathcal{D}} + \beta\bar{w}_0 \otimes \bar{w}_0)\bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}})}{(\mathbf{S}_{00}\bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}})} \stackrel{(49)}{\geq} \frac{((\mathbf{B}_{\mathcal{D}} + \beta\bar{w}_0 \otimes \bar{w}_0)\bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}})}{(\hat{\mathbf{S}}_{00}\bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}})} \\ & \stackrel{(55),(54)}{\geq} C \frac{((\mathbf{B}_{\mathcal{D}} + \frac{1}{R^2}\mathbf{M}_{\mathcal{D}})\bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}})}{(\hat{\mathbf{A}}\hat{u}, \hat{u})} \stackrel{(56)}{\geq} \frac{C\|u_h^{\mathcal{D}}\|_{\mathcal{D}_h}^2}{(C^* + 1)\|u_h^{\mathcal{D}}\|_{\mathcal{D}_h}^2} = \frac{C}{(C^* + 1)} =: \mu_{\star}, \quad \text{with } \beta = \frac{h^2}{R^2} \end{aligned}$$

where μ_{\star} is independent of $h > 0$.

From the obtained above bounds, we have

$$\mu_{\star} \leq \frac{(\mathbf{B}_{\mathcal{D}}\bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}})}{(\mathbf{S}_{00}\bar{u}_{\mathcal{D}}, \bar{u}_{\mathcal{D}})} \leq 1, \quad \text{for } \mathbb{R}^n \ni \bar{u}_{\mathcal{D}} \perp \ker \mathbf{B}_{\mathcal{D}}.$$

□

3.4 Notes on Lanczos algorithm with the block-diagonal preconditioner \mathcal{P}

The preconditioned Lanczos procedure of minimized iterations can be used for solving algebraic systems with symmetric and positive semidefinite matrices. In this section, we propose a preconditioner for solving (21).

The theoretical justification of the usage of a preconditioner (26) where the blocks \mathcal{P}_A and \mathcal{P}_B are spectrally equivalent to \mathbf{A} and $\mathbf{B}\mathbf{A}^{-1}\mathbf{B}^T$, respectively, was shown in [10]. With provided above theoretical considerations, in our practical implementation of the generalized Lanczos method of minimized iterations, we use the following block-diagonal preconditioner:

$$\mathcal{P} = \begin{bmatrix} \mathcal{P}_A & 0 \\ 0 & \mathbf{B}_D \end{bmatrix}, \quad (57)$$

where one can choose any typical preconditioner \mathcal{P}_A for the symmetric and positive-definite matrix \mathbf{A} . This in particular might be \mathbf{A} itself as we use it below. Define

$$\mathcal{H} = \mathcal{P}^\dagger = \begin{bmatrix} \mathcal{P}_A^{-1} & 0 \\ 0 & [\mathbf{B}_D]^\dagger \end{bmatrix}, \quad (58)$$

and a new scalar product

$$(\bar{x}, \bar{y})_{\mathcal{H}} := (\mathcal{H}\bar{x}, \bar{y}), \quad \text{for all } \bar{x}, \bar{y} \in \mathbb{R}^{N+n},$$

and consider the preconditioned Lanczos iterations $\bar{z}^k = \begin{bmatrix} \bar{u}^k \\ \bar{\lambda}^k \end{bmatrix} \in \mathbb{R}^{N+n}$, $k \geq 1$:

$$\bar{z}^k = \bar{z}^{k-1} - \beta_k \bar{y}_k,$$

where

$$\beta_k = \frac{(\mathcal{A}_\varepsilon \bar{z}^{k-1} - \bar{\mathcal{F}}, \mathcal{A}_\varepsilon \bar{y}_k)_{\mathcal{H}}}{(\mathcal{A}_\varepsilon \bar{y}_k, \mathcal{A}_\varepsilon \bar{y}_k)_{\mathcal{H}}}.$$

and

$$y_k = \begin{cases} \mathcal{H}(\mathcal{A}_\varepsilon \bar{z}^0 - \bar{\mathcal{F}}), & k = 1 \\ \mathcal{H}\mathcal{A}_\varepsilon \bar{y}_1 - \alpha_2 \bar{y}_1, & k = 2 \\ \mathcal{H}\mathcal{A}_\varepsilon \bar{y}_{k-1} - \alpha_k \bar{y}_{k-1} - \gamma_k \bar{y}_{k-2}, & k > 2, \end{cases}$$

with

$$\alpha_k = \frac{(\mathcal{A}_\varepsilon \mathcal{H}\mathcal{A}_\varepsilon \bar{y}_{k-1}, \mathcal{A}_\varepsilon \bar{y}_{k-1})_{\mathcal{H}}}{(\mathcal{A}_\varepsilon \bar{y}_{k-1}, \mathcal{A}_\varepsilon \bar{y}_{k-1})_{\mathcal{H}}}, \quad \gamma_k = \frac{(\mathcal{A}_\varepsilon \mathcal{H}\mathcal{A}_\varepsilon \bar{y}_{k-1}, \mathcal{A}_\varepsilon \bar{y}_{k-1})_{\mathcal{H}}}{(\mathcal{A}_\varepsilon \bar{y}_{k-2}, \mathcal{A}_\varepsilon \bar{y}_{k-2})_{\mathcal{H}}}.$$

Here, we recall that the matrix \mathbf{B}_D is singular, however, as evident from the algorithm above one actually never needs to use its pseudo-inverse at all. Indeed, this is due to the block-diagonal structure of \mathcal{H} (58), and block form of the original matrix \mathcal{A}_ε (16)-(17).

4 Numerical Results

In this section, we use three examples to show the numerical advantages of the Lanczos iterative scheme with the preconditioner \mathcal{P} defined in (57) over the existing preconditioned conjugate gradient method.

Our numerical experiments are performed by implementing the described above Lanczos algorithm for the problem (4)-(5), where the domain Ω is chosen to be a disk of radius 5 with $m = 37$ identical circular inclusions \mathcal{D}^i , $i \in \{1, \dots, m\}$. Inclusions are equally spaced. The function f of the right hand side of (4) is chosen to be a constant, $f = 50$.

In the first set of experiments the values of ε of (5) are going to be identical in all inclusions and vary from 10^{-1} to 10^{-8} . In the second set of experiments we consider four groups of particles with

the same values of ε in each group that vary from 10^{-4} to 10^{-7} . In the third set of experiments we decrease the distance between neighboring inclusions.

The initial guess z^0 is a random vector that was fixed for all experiments. The stopping criteria is the Euclidian norm of the relative residual $(\mathcal{A}_\varepsilon z^k - \bar{\mathcal{F}})/\bar{\mathcal{F}}$ being less than a fixed tolerance constant.

We test our results against standard `pcg` function of MATLAB[®] with $\mathcal{P}_A = \mathbf{A}$. The same matrix is also used in the implementation of the described above Lanczos algorithm. In the following tables **PCG** stands for preconditioned conjugate gradient and **PL** stands for preconditioned Lanczos.

Experiment 1. For the first set of experiments we consider particles \mathcal{D}^i of radius $R = 0.45$ in the disk Ω . This choice makes distance d between neighboring inclusions approximately equal to the radii of inclusions. The triangular mesh Ω_h has $N = 32,567$ nodes. Tolerance is chosen to be equal to 10^{-4} . This experiment concerns the described problem with parameter ε being the same in each inclusion. Table 1 shows the number of iterations corresponding to the different values of ε .

Table 1: Iteration numbers depending on values of ε , $N = 32,567$

| ε | 10^{-1} | 10^{-2} | 10^{-3} | 10^{-4} | 10^{-5} | 10^{-6} | 10^{-7} | 10^{-8} |
|---------------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| PCG | 10 | 20 | 32 | 40 | 56 | 183 | 302 | 776 |
| PL | 33 | 37 | 37 | 37 | 37 | 37 | 37 | 37 |

Based on these results, we first observe that our **PL** method requires less iterations as ε goes less than 10^{-4} . We also notice that number of iterations in the Lanczos algorithm does not depend on ε .

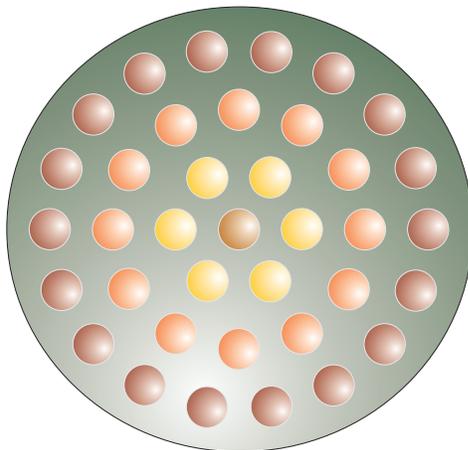


Figure 3: The domain Ω with highly conducting inclusions \mathcal{D}^i of four groups

Experiment 2. In this experiment we leave radii of the inclusions to be the same, namely, $R = 0.45$. Tolerance is chosen to be 10^{-6} . The key point of this set of experiments is to split inclusions in four groups and assign different values of ε to each group. The first group consists of one inclusion – in the center – with the coefficient $\varepsilon = \varepsilon_1$, whereas the second, third and fourth groups are comprised of the disks in the second, third and fourth circular layers of inclusions with coefficients ε_2 , ε_3 and ε_4 respectively, see Fig. 3. We perform this type of experiments for three different triangular meshes with the total number of nodes $N = 5,249$, $N = 12,189$ and

$N = 32,567$. Tables 2, 3, and 4 below show the number of iterations corresponding to three meshes respectively.

Table 2: Iteration numbers for values of ε different in each group of inclusions, $N = 5,249$

| ε_1 | ε_2 | ε_3 | ε_4 | PCG | PL |
|-----------------|-----------------|-----------------|-----------------|------------|-----------|
| 10^{-5} | 10^{-5} | 10^{-4} | 10^{-4} | 217 | 39 |
| 10^{-5} | 10^{-5} | 10^{-4} | 10^{-3} | 208 | 39 |
| 10^{-6} | 10^{-5} | 10^{-4} | 10^{-3} | 716 | 39 |
| 10^{-7} | 10^{-6} | 10^{-5} | 10^{-4} | 571 | 39 |

Table 3: Iteration numbers for values of ε different in each group of inclusions, $N = 12,189$

| ε_1 | ε_2 | ε_3 | ε_4 | PCG | PL |
|-----------------|-----------------|-----------------|-----------------|------------|-----------|
| 10^{-5} | 10^{-5} | 10^{-4} | 10^{-4} | 116 | 39 |
| 10^{-5} | 10^{-5} | 10^{-4} | 10^{-3} | 208 | 39 |
| 10^{-6} | 10^{-5} | 10^{-4} | 10^{-3} | 457 | 39 |
| 10^{-7} | 10^{-6} | 10^{-5} | 10^{-4} | 454 | 39 |

Table 4: Iteration numbers for values of ε different in each group of inclusions, $N = 32,567$

| ε_1 | ε_2 | ε_3 | ε_4 | PCG | PL |
|-----------------|-----------------|-----------------|-----------------|------------|-----------|
| 10^{-5} | 10^{-5} | 10^{-4} | 10^{-4} | 311 | 35 |
| 10^{-5} | 10^{-5} | 10^{-4} | 10^{-3} | 311 | 35 |
| 10^{-6} | 10^{-5} | 10^{-4} | 10^{-3} | 697 | 35 |
| 10^{-7} | 10^{-6} | 10^{-5} | 10^{-4} | 693 | 35 |

These results yield that **PL** requires much less iterations than **PCG** with the number of iterations still being independent of both the contrast ε and the mesh size h .

Experiment 3. Next we take the same set up of 37 inclusions and decrease the distance between them by making radius of each inclusion larger, now let $R = 0.56$. Radius of each inclusion is now twice larger than the distance d . Tolerance is chosen to be 10^{-6} . The triangular mesh Ω_h has $N = 6,329$ nodes. Table 5 shows the number of iterations in this case.

We notice that number of iterations goes up for both **PCG** and **PL**, while this number remaining independent of ε for **PL**.

Further we continue to decrease the distance d , taking $R = 0.62$. This number makes radius of each inclusion four times larger than the distance between two neighboring inclusions. Tolerance is fixed at 10^{-6} . The triangular mesh Ω_h has $N = 6,699$ nodes. We observed that **PL** method did not reach the desired tolerance in 1,128 iterations. The result is expected and shows that assumption about the ration between R and d is crucial for existence of harmonic extension.

Table 5: Iteration numbers for values of ε different in each group of inclusions, $N = 6,329$

| ε_1 | ε_2 | ε_3 | ε_4 | PCG | PL |
|-----------------|-----------------|-----------------|-----------------|-----|----|
| 10^{-5} | 10^{-5} | 10^{-4} | 10^{-4} | 799 | 61 |
| 10^{-7} | 10^{-6} | 10^{-5} | 10^{-4} | 859 | 61 |

5 Conclusions

This paper focuses a construction of the robust preconditioner (57) for the Lancsoz iterative scheme that can be used in order to solve high-contrast PDEs of the type (4)-(5). A typical FEM discretization yields an ill-conditioning matrix when the contrast in σ becomes high (i.e. $\varepsilon \ll 1$). We propose a saddle point formulation of the given problem with the symmetric indefinite matrix and consequently construct the corresponding preconditioner that yields a robust numerical approximation of (4)-(5). The main feature of this novel and elegant approach is that we precondition the given linear system with a *symmetric indefinite matrix*. Our numerical results have shown the effectiveness of the proposed preconditioner for these type of problems.

6 Appendix

Here we prove Lemma 1.

Proof. Without loss of generality, here we also assume that all $\varepsilon_i = \varepsilon$, $i \in \{1, \dots, m\}$. Hereafter we denote by C a positive constant that is independent of ε .

Subtract first equations of (8) and (22) and multiply by $\bar{u}_\varepsilon - \bar{u}_0$ to obtain

$$(\mathbf{A}(\bar{u}_\varepsilon - \bar{u}_0), \bar{u}_\varepsilon - \bar{u}_0) + (\mathbf{B}^T(\bar{\lambda}_\varepsilon - \bar{\lambda}_0), \bar{u}_\varepsilon - \bar{u}_0) = \bar{0}.$$

Recall, that the matrix \mathbf{A} is SPD then

$$(\mathbf{A}\xi, \xi) \geq \mu_1(\mathbf{A})\|\xi\|^2, \quad \forall \xi \in \mathbb{R}^N,$$

where $\mu_1(\mathbf{A}) > 0$ is the *minimal eigenvalue* of \mathbf{A} .

Making use of the second equation of (15) we have

$$\mu_1(\mathbf{A})\|\bar{u}_\varepsilon - \bar{u}_0\|^2 \leq -(\varepsilon \mathbf{B}_D \bar{\lambda}_\varepsilon, \bar{\lambda}_\varepsilon) + (\varepsilon \mathbf{B}_D \bar{\lambda}_\varepsilon, \bar{\lambda}_0) \leq (\varepsilon \mathbf{B}_D \bar{\lambda}_\varepsilon, \bar{\lambda}_0),$$

where we used the fact that \mathbf{B}_D is positive semidefinite. Then

$$\|\bar{u}_\varepsilon - \bar{u}_0\|^2 \leq \varepsilon \|\mathbf{B}_D \bar{\lambda}_\varepsilon\|. \quad (59)$$

To estimate $\mathbf{B}_D \bar{\lambda}_\varepsilon$ we multiply the first equation of (15) by $\mathbf{A} \mathbf{B}^T \bar{\lambda}_\varepsilon$:

$$(\mathbf{A} \bar{u}_\varepsilon, \mathbf{A} \mathbf{B}^T \bar{\lambda}_\varepsilon) + (\mathbf{B}^T \bar{\lambda}_\varepsilon, \mathbf{A} \mathbf{B}^T \bar{\lambda}_\varepsilon) = (\bar{\mathbf{F}}, \mathbf{A} \mathbf{B}^T \bar{\lambda}_\varepsilon),$$

that yields

$$\mu_1(\mathbf{A})\|\mathbf{B}^T \bar{\lambda}_\varepsilon\|^2 \leq C\|\bar{\mathbf{F}} - \mathbf{A} \bar{u}_\varepsilon\| \|\mathbf{B}^T \bar{\lambda}_\varepsilon\|.$$

Note that $\mathbf{B}^T \bar{\lambda}_\varepsilon = \mathbf{B}_D \bar{\lambda}_\varepsilon$, hence,

$$\|\mathbf{B}_D \bar{\lambda}_\varepsilon\| \leq C\|\bar{\mathbf{F}} - \mathbf{A} \bar{u}_\varepsilon\|, \quad (60)$$

so collecting estimates (59)-(60), it remains to show $\|\bar{u}_\varepsilon\|$ is bounded. For that we multiply the first equation of (8) by \bar{u}_ε and obtain

$$(\mathbf{A}\bar{u}_\varepsilon, \bar{u}_\varepsilon) + (\mathbf{B}^T \bar{\lambda}_\varepsilon, \bar{u}_\varepsilon) = (\bar{\mathbf{F}}, \bar{u}_\varepsilon),$$

that yields

$$\|\bar{u}_\varepsilon\| \leq C.$$

□

References

- [1] B. Aksoylu, I. G. Graham, H. Klie, and R. Scheichl, “Towards a rigorously justified algebraic preconditioner for high-contrast diffusion problems”, *Computing and Visualization in Science*, **11:4-6**, 2008, pp. 319–331
- [2] J. Aarnes, and T. Y. Hou, “Multiscale domain decomposition methods for elliptic problems with high aspect ratios”, *Acta Mathematicae Applicatae Sinica. English Series*, **18:1**, 2002, pp. 63–76
- [3] O. Axelsson, *Iterative Solution Methods*, Cambridge University Press, 1994
- [4] Z.-Z. Bai, M. K. Ng, and Z.-Q. Wang, “Constraint preconditioners for symmetric indefinite matrices”, *SIAM Journal on Matrix Analysis and Applications*, **31:2**, 2009, pp. 410–433
- [5] M. Benzi, G. H. Golub, and J. Liesen, “Numerical solution of saddle point problems”, *Acta Numerica*, **14**, 2005, pp. 1–137
- [6] M. Benzi, and A. J. Wathen, “Some preconditioning techniques for saddle point problems”, in *Model order reduction: theory, research aspects and applications*, Springer, Berlin, **13**, 2008, pp. 195–211
- [7] V. M. Calo, Y. Efendiev, and J. Galvis, “Asymptotic expansions for high-contrast elliptic equations”, *Mathematical Models and Methods in Applied Sciences*, **24:3**, 2014, pp. 465–494
- [8] H. C. Elman, D. J. Silvester, and A. J. Wathen, *Finite elements and fast iterative solvers: with applications in incompressible fluid dynamics*, in *Numerical Mathematics and Scientific Computation*, Oxford University Press, New York, 2005
- [9] R. Glowinski, and Yu. Kuznetsov, “On the solution of the Dirichlet problem for linear elliptic operators by a distributed Lagrange multiplier method”, *Comptes Rendus de l’Académie des Sciences. Série I. Mathématique*, **327:7**, 1998, pp. 693–698
- [10] Yu. Iliash, T. Rossi, and J. Toivanen, “Two iterative methods to solve the Stokes problem”, *Technical Report No. 2. Lab. Sci. Comp.*, Dept. Mathematics, University of Jyväskylä. Jyväskylä, Finland, 1993.
- [11] C. Keller, N. I. M. Gould, and A. J. Wathen, “Constraint preconditioning for indefinite linear systems”, *SIAM Journal on Matrix Analysis and Applications*, **21:4**, 2000, pp. 1300–1317
- [12] Yu. Kuznetsov, “Efficient iterative solvers for elliptic finite element problems on nonmatching grids”, *Russian Journal of Numerical Analysis and Mathematical Modelling*, **10:3**, 1995, pp. 187–211

- [13] Yu. Kuznetsov, “Preconditioned iterative methods for algebraic saddle-point problems”, *Journal of Numerical Mathematics*, **17:1**, 2009, pp. 67-75
- [14] Yu. Kuznetsov, and G. Marchuk, “Iterative methods and quadratic functionals”. In *Méthodes de l’Informatique-4*, eds. J.-L. Lions and G. Marchuk, pp. 3–132, Paris, 1974 (In French)
- [15] L. Lukšan, and J. Vlček, “Indefinitely preconditioned inexact Newton method for large sparse equality constrained non-linear programming problems”, *Numerical Linear Algebra with Applications*, **5:3**, 1998, pp. 219–247
- [16] C. C. Paige, “Computational variants of the Lanczos method for the eigenproblem”, *Journal of the Institute of Mathematics and its Applications*, **10**, 1972, pp. 373–381
- [17] A. Toselli, and O. Widlund, “Domain decomposition methods – algorithms and theory”, Springer Series in Computational Mathematics, **34**, Springer-Verlag, Berlin, 2005
- [18] E. L. Wachspress, “Iterative solution of elliptic systems, and applications to the neutron diffusion equations of reactor physics”, Prentice-Hall, Inc., Englewood Cliffs, N.J., 1966
- [19] A. J. Wathen, “Preconditioning”, *Acta Numerica*, **24**, 2015, pp. 329–376
- [20] O. B. Widlund, “An Extension Theorem for Finite Element Spaces with Three Applications”, Chapter *Numerical Techniques in Continuum Mechanics* in *Notes on Numerical Fluid Mechanics*, **16**, 1987, pp. 110–122
- [21] X. Wu, B. P. B. Silva, and J. Y. Yuan, “Conjugate gradient method for rank deficient saddle point problems”, *Numerical Algorithms*, **35:2–4**, 2004, pp. 139–154