# On the Moore-Penrose inverse in solving saddle-point systems with singular diagonal blocks

R. Kučera[1,*], T. Kozubek[1], A. Markopoulos[1], and J. Machalová[2]

[1] *VŠB-Technical University Ostrava, Czech Republic*
[2] *Palacký University Olomouc, Czech Republic*

## SUMMARY

The paper deals with the role of the generalized inverses in solving saddle-point systems arising naturally in the solution of many scientific and engineering problems when FETI based domain decomposition methods are used to their numerical solution. It is shown that the Moore-Penrose inverse may be obtained in this case at negligible cost by projecting an arbitrary generalized inverse using orthogonal projectors. Applying an eigenvalue analysis based on the Moore-Penrose inverse, it is proved for simple model problems that the number of conjugate gradient iterations required for the solution of associate dual systems does not depend on discretization norms. The theoretical results are confirmed by numerical experiments with linear elasticity problems. Copyright © 2011 John Wiley & Sons, Ltd.

KEY WORDS:    Moore-Penrose inverse; orthogonal projectors; saddle-point systems; domain decomposition methods; condition number

*The first and the fourth authors would like to dedicate this paper to their former teacher of numerical mathematics doc. Jiří Kobza on the occasion of his 70th birthday.*

## 1. INTRODUCTION

The paper deals with solving large *saddle-point systems* with singular diagonal blocks (see [2] and references therein), i.e., solving the problem to find $(\bar{u}, \bar{\lambda}) \in \mathbb{R}^n \times \mathbb{R}^m$ satisfying:

$$\left( \begin{array}{cc} A & B^\top \\ B & 0 \end{array} \right) \left( \begin{array}{c} u \\ \lambda \end{array} \right) = \left( \begin{array}{c} f \\ g \end{array} \right), \tag{1}$$

where $A \in \mathbb{R}^{n \times n}$ is a symmetric, positive semi-definite matrix, $B \in \mathbb{R}^{m \times n}$ is a full-row rank matrix, $f \in \mathbb{R}^n$, $g \in \mathbb{R}^m$, and $m \ll n$. Such systems arise in many scientific and

engineering applications when FETI (Finite Element Tearing and Interconnecting) based domain decomposition methods [12, 9] or fictitious domain methods [17, 16] are used to the numerical solution. Here, we will pay our attention to a variant of the FETI method called Total FETI [9].

The FETI methods belong to the most efficient domain decomposition techniques for the numerical solution of boundary value problems described by elliptic partial differential equations (PDEs). In this context, $A$ plays the role of the stiffness matrix, $B$ is the "gluing" matrix, and $\bar{\lambda}$ is the vector of *Lagrange multipliers* enforcing the continuity of the PDE solution. There are two main benefits of the FETI approach. Firstly, the stiffness matrix has a block diagonal structure that enables us to handle the blocks in parallel. Secondly, the condition number of all blocks and, consequently of the whole stiffness matrix, may be independent on the size of the discrete problem (under additional assumptions on the finite element partitions). It is well-known that convergence of conjugate gradient type methods is determined by the condition number of the system matrix [14]. Therefore the number of iterations we need to get a solution to (1) with a given accuracy may be independent on the size of the discrete problem, as well. This property is known as the *scalability* of the method [8].

The classical FETI algorithm [12] consists in eliminating the first unknown $u$ from (1) and solving the resulting dual system in terms of $\lambda$ iteratively. As the diagonal blocks of $A$ may be singular, the elimination requires to use a generalized inverse to $A$ and a basis of the kernel-space of $A$. One of the reasons to develop other variants of the FETI method was the effort to overcome difficulties with computing the action of generalized inverses and identifying kernel-spaces. The FETI-DP variant [11, 19] modifies the original FETI method so that $A$ is nonsingular. Then the kernel-space is trivial and the inverse to $A$ exists. The opposite strategy gave rise to the TFETI (Total FETI) method [9] in which also the Dirichlet boundary conditions of the PDE problem are enforced by Lagrange multipliers. In this case, the kernel-space is as large as possible, since all diagonal blocks of $A$ are subdomain stiffness matrices to the original PDEs with the *pure* homogeneous Neumann boundary conditions. The advantage is that the kernel-space basis is known à-priori and, in addition, it may be assembled by mechanical arguments at negligible cost.

The main goal of our paper is to show how to handle the generalized inverse to $A$ in the TFETI method. As stiffness matrices to elliptic PDE problems are symmetric, positive semi-definite, the generalized Cholesky factorization can be applied [14]. Then, by inverting the regular part of the Cholesky factor, one can easily get a generalized inverse to $A$. Unfortunately, the factorization procedure is sensitive to round-off errors, since the zero pivots must be recognized. Therefore Farhat and Gèradin [10] proposed to replace the Cholesky factorization by the singular value decomposition (SVD) immediately when a suspected zero pivot appears. This technique leads to the robust algorithm but still dependent on the tolerance for zero pivot detection. The method was further developed by Fragakis and Papadrakakis [22] who proposed an effective mechanism for identification of zero pivots. Another development of this idea was done by Dostál et. al. [6] who proposed to modify the Farhat and Géradin procedure in order to eliminate the identification problem and to reduce the factorization to an à-priori defined well-conditioned positive definite diagonal block of $A$.

Let us note that the Moore-Penrose (MP) inverse is generally the best generalized inverse for the iterative solution of problems with singular matrices. The reason is that it minimizes (among all generalized inverses) the norm of the computed vector [7] that keeps the used iterative method as stable as it is possible. This fact is highly important when the corresponding

saddle-point systems (1) are large and ill-conditioned. Our results will show that the (exact) MP inverse can be obtained from an arbitrary (stable computable) generalized inverse by its projection using the orthogonal projector on the image-space of $A$. Since the orthogonal projector is available in the TFETI method due to the knowledge of the kernel-space basis, the MP inverse may be easily implemented. Our idea is closely related to Pyle's algorithm [25] which, however, assembles the MP inverse projecting the generalized inverse computed by the orthogonal factorization.

The paper is organized as follows. Section 2 deals with a generalized inverse to an arbitrary (rectangular) matrix. We derive the three-condition characterization of the MP inverse based on identifying its kernel-space and image-space. Then we prove how to get the MP inverse from an arbitrary generalized inverse using orthogonal projectors. In Section 3, we discuss the use of the MP inverse in solving saddle-point systems. The eigenvalue analysis shows that the condition number of the corresponding dual system is bounded by the condition number of $A$ and $BB^\top$. Applying this result in Section 4 for simple model PDE problems in one and two space dimensions (1D and 2D) we prove that the number of conjugate gradient iterations required for solving the dual system does not depend on its size. Finally in Section 5, we demonstrate the theoretical results on the numerical solution of more complex problems arising from linear elasticity in three space dimensions (3D).

## 2. GENERALIZED INVERSES AND PROJECTORS

In this section, we prove three conditions determining the Moore-Penrose (MP) inverse to a rectangular matrix. Then we show how to transform an arbitrary generalized inverse to the MP one using orthogonal projectors. Before getting these results we start with preliminaries.

Let $\mathbb{R}^{m \times n}$ be the set of $m \times n$ real matrices and let $A \in \mathbb{R}^{m \times n}$. The symbols $Ker\,A$ and $Im\,A$ stand for the kernel (kernel-space) and the image (image-space) of $A$, respectively. The rank of $A$ is defined by $r(A) := \dim Im\,A$ and obviously $r(A) = r(A^\top)$, where $A^\top$ is the transpose to $A$. Finally, let $I$ and $0$ denote the identity and zero matrices, respectively.

By a *generalized inverse* to $A$ we call such $X \in \mathbb{R}^{n \times m}$ that satisfies the following equation:

$$A = AXA. \tag{2}$$

Let us note that there is a generalized inverse for any $A$ but it is not uniquely determined by (2). On the other hand the MP inverse to $A$, denoted here by $A^\dagger$, is a particular generalized inverse that is uniquely defined for any $A$. There are various definitions of $A^\dagger$. Here we will define $A^\dagger$ by the singular value decomposition (SVD) of $A$.

**Theorem 2.1.** *Let $A \in \mathbb{R}^{m \times n}$ be of rank $r = r(A)$. There are orthogonal matrices $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$ and diagonal matrix $\Sigma \in \mathbb{R}^{m \times n}$ with non-negative entries so that*

$$A = U\Sigma V^\top. \tag{3}$$

*Moreover, the diagonal entries of $\Sigma$ can be sorted in the decreasing order so that*

$$\Sigma = \left( \begin{array}{cc} \widehat{\Sigma} & 0 \\ 0 & 0 \end{array} \right), \tag{4}$$

*where $\widehat{\Sigma} \in \mathbb{R}^{r \times r}$ is the nonsingular part of $\Sigma$.*

http://hdl.handle.net/10084/94957 09/08/2012

*Proof.* See [1].                                                                                                                □

By the SVD of $A$ we understand $U$, $V$, and $\Sigma$ satisfying (3) and (4). Let us note that the SVD is uniquely determined by (3) and (4) for any $A$. In addition, we split $U$ and $V$ accordingly to (4), i.e., $U_1 \in \mathbb{R}^{m \times r}$, $U_2 \in \mathbb{R}^{m \times m-r}$, $V_1 \in \mathbb{R}^{n \times r}$, and $V_2 \in \mathbb{R}^{n \times n-r}$ are such that $U$, $V$ take the form $U = (U_1, U_2)$, $V = (V_1, V_2)$, respectively. Then

$$A = (U_1, U_2) \begin{pmatrix} \widehat{\Sigma} & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} V_1^\top \\ V_2^\top \end{pmatrix} \tag{5}$$

and

$$A^\top = (V_1, V_2) \begin{pmatrix} \widehat{\Sigma} & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} U_1^\top \\ U_2^\top \end{pmatrix}. \tag{6}$$

The following lemma shows the orthogonal decompositions of $\mathbb{R}^m$ and $\mathbb{R}^n$ defined by $A$ and $A^\top$, respectively.

**Lemma 2.1.** *It holds:*

$$\mathbb{R}^m = Im\, A \oplus Ker\, A^\top \quad and \quad Im\, A \perp Ker\, A^\top; \tag{7}$$
$$\mathbb{R}^n = Im\, A^\top \oplus Ker\, A \quad and \quad Im\, A^\top \perp Ker\, A. \tag{8}$$

*Proof.* The columns of $U$ from the SVD of $A$ are basis in $\mathbb{R}^m$. Further, (5) and (6) imply that the columns of $U_1$ and $U_2$ are the orthogonal bases in $Im\, A$ and $Ker\, A^\top$, respectively. The relations (7) follow from $U = (U_1, U_2)$. The proof of (8) is analogous.                □

Now we give the definition of the MP inverse based on the SVD.

**Definition 2.1.** *Let $A \in \mathbb{R}^{m \times n}$ be given. By the MP inverse to $A$ we call the matrix $A^\dagger$ defined by*

$$A^\dagger := (V_1, V_2) \begin{pmatrix} \widehat{\Sigma}^{-1} & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} U_1^\top \\ U_2^\top \end{pmatrix}, \tag{9}$$

*where $U$, $V$, and $\widehat{\Sigma}$ are determined uniquely by the SVD of $A$.*

It is easy to verify that the MP inverse given by (9) is the generalized inverse, i.e., (2) is satisfied for $X = A^\dagger$. The following lemma holds immediately.

**Lemma 2.2.** *Let $A^\dagger$ be the MP inverse to $A$. Then*

$$Ker\, A^\dagger = Ker\, A^\top, \quad Im\, A^\dagger = Im\, A^\top. \tag{10}$$

*Proof.* Compare (9) and (6).                                                                                □

Let us note that Moore's definition [20, 15] of the MP inverse, more or less forgotten, characterizes the MP inverse $A^\dagger$ by the following three conditions:

$$A = AA^\dagger A, \quad Im\, A^\dagger \subseteq Im\, A^\top, \quad Im\, (A^\dagger)^\top \subseteq Im\, A.$$

In spite of Moore's definition, we will prove that the MP-inverse is the generalized inverse satisfying (10). We will need several auxiliary results. First of all we recall the well-known conditions used by Penrose [23] to define the MP inverse.

**Lemma 2.3.** *Let $A \in \mathbb{R}^{m \times n}$ be given. Then $X$ is the MP inverse to $A$, i.e. $X = A^\dagger$, iff*

$$A = AXA, \quad XA = (XA)^\top, \quad AX = (AX)^\top, \quad X = XAX. \tag{11}$$

*Proof.* There is a unique $X$ determined by (11); see [1]. One can verify that $X = A^\dagger$ given by (9) satisfies all equations (11). $\qquad\square$

**Lemma 2.4.** *Let $P$ be a square matrix. If $Im\,P \perp Im\,(I - P)$, then $P$ is symmetric.*

*Proof.* The orthogonality of $Im\,P$ and $Im\,(I - P)$ is equivalent to $P^\top(I - P) = 0$ so that $P^\top = P^\top P = (P^\top P)^\top = (P^\top)^\top = P$. $\qquad\square$

Let us recall that any square matrix $P$ satisfying $P^2 = P$ is called the *projector* on $Im\,P$. Then $P^\top$ is the projector on $Im\,P^\top$, since $(P^\top)^2 = P^\top$. Moreover, if $Im\,P \perp Im\,(I - P)$, the projector $P$ is called *orthogonal*.

**Lemma 2.5.** *Let $P$ be a projector. Then $P$ is orthogonal iff it is symmetric.*

*Proof:* As Lemma 2.4 holds, it remains to prove that the symmetry implies the orthogonality. It is $P^\top(I - P) = P^\top - P^\top P = P - P^2 = 0$. $\qquad\square$

**Lemma 2.6.** *Let $X$ be a generalized inverse to $A$. Then*
*(i) $Y := AX$ is the projector on $Im\,A$;*
*(ii) $I - Y^\top$ is the projector on $Ker\,A^\top$;*
*(iii) $Z := (XA)^\top$ is the projector on $Im\,A^\top$;*
*(iv) $I - Z^\top$ is the projector on $Ker\,A$.*

*Proof:* As $Y^2 = AXAX = AX = Y$, we see that $Y$ is the projector on $Im\,Y \subseteq Im\,A$. For $x \in Im\,A$, $x = Ay$, we obtain $Yx = AXAy = Ay = x$ so that $Im\,Y = Im\,A$ and (i) holds. Further, $(I - Y^\top)^2 = I - 2Y^\top + (Y^\top)^2 = I - Y^\top$ implies that $I - Y^\top$ is the projector on $Im\,(I - Y^\top)$. Moreover, $A^\top(I - Y^\top) = A^\top - A^\top X^\top A^\top = 0$ yields $Im\,(I - Y^\top) \subseteq Ker\,A^\top$. For $x \in Ker\,A^\top$, we have $(I - Y^\top)x = x - X^\top A^\top x = x$ so that $Im\,(I - Y^\top) = Ker\,A^\top$ and therefore (ii) is valid. The proof of (iii) and (iv) is analogous. $\qquad\square$

**Corollary 2.1.** *Let $X$ be a generalized inverse to $A$. It holds:*
*(a) $Im\,AX = Im\,A$;*
*(b) $Im\,(I - (AX)^\top) = Ker\,A^\top$;*
*(c) $Im\,(XA)^\top = Im\,A^\top$;*
*(d) $Im\,(I - XA) = Ker\,A$;*
*(e) $Ker\,(AX)^\top = Ker\,A^\top$;*
*(f) $Ker\,(I - AX) = Im\,A$;*
*(g) $Ker\,XA = Ker\,A$;*
*(h) $Ker\,(I - (XA)^\top) = Im\,A^\top$.*

*Proof.* The statements (a)-(d) follow from Lemma 2.6 and (e)-(h) are the equalities of the respective orthogonal complements. $\qquad\square$

Now we are able to prove the three conditions determining the MP inverse.

**Theorem 2.2.** *Let $A \in \mathbb{R}^{m \times n}$ be given. Then $X$ is the MP inverse to $A$, i.e. $X = A^\dagger$, iff*
*(i) $A = AXA$,*
*(ii) $Im\, X = Im\, A^\top$,*
*(iii) $Ker\, X = Ker\, A^\top$.*

*Proof.* As Lemma 2.2 holds, it remains to prove the converse implication "$\Leftarrow$". We will verify $(11)_1$-$(11)_4$. The first condition $(11)_1$ holds since it is (i). The trivial inclusion $Im\, X \supseteq Im\, XA$ and (ii) imply $Im\, A^\top \supseteq Im\, XA$. By Corollary 2.1(c) we obtain $r(A^\top) = r((XA)^\top) = r(XA)$ so that $Im\, A^\top = Im\, XA$. Using (8) and Corollary 2.1(d) we arrive at

$$Im\, XA = Im\, A^\top \perp Ker\, A = Im\, (I - XA).$$

Therefore Lemma 2.4 implies $XA = (XA)^\top$ that is $(11)_2$. To prove $(11)_3$ we start with (iii) and Corollary 2.1(e) that give $Ker\, X = Ker\, (AX)^\top$. Passing to the orthogonal complements we get $Im\, X^\top = Im\, AX$ that yields $r(X^\top) = r(AX) = r((AX)^\top)$. This equality together with the obvious inclusion $Im\, X^\top \supseteq Im\, (AX)^\top$ lead to $Im\, X^\top = Im\, (AX)^\top$. By (8), again (iii), and Corollary 2.1(b) we get

$$Im\, (AX)^\top = Im\, X^\top \perp Ker\, X = Ker\, A^\top = Im\, (I - (AX)^\top).$$

Now Lemma 2.4 implies $(AX)^\top = AX$ that is $(11)_3$. Moreover, we obtain

$$Im\, X^\top \perp Im\, (I - AX)$$

or, equivalently, $X(I - AX) = 0$ that proves $(11)_4$. □

**Example 2.1.** Let us consider the (symmetric) matrix

$$A = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$

with $Im\, A$ and $Ker\, A$ generated by the vectors $(1,1)^\top$ and $(1,-1)^\top$, respectively; see Figure 1. Let us define $S_X := \sum_{ij} x_{ij}$ for any $X = (x_{ij}) \in \mathbb{R}^{2 \times 2}$. It is readily seen that

$$AXA = \begin{pmatrix} S_X & S_X \\ S_X & S_X \end{pmatrix} = A$$

holds iff $S_X = 1$. Therefore generalized inverses to $A$ can be obtained from every $M \in \mathbb{R}^{2 \times 2}$ such that $S_M \neq 0$ by

$$X := S_M^{-1} M. \tag{12}$$

Next, we will show how to transform an arbitrary generalized inverse $X$ to the MP one. Let us introduce $M$ with the image and the kernel generated by the nonzero vectors $v = (v_1, v_2)^\top$ and $w = (w_1, w_2)^\top$, respectively, i.e.,

$$M := \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} (w_2, -w_1).$$

Let us define the generalized inverse $X$ by (12); see Figure 2. Note that

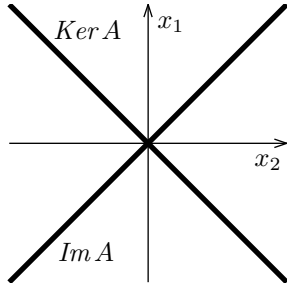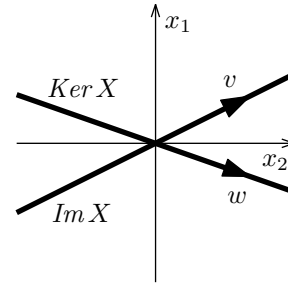$$S_M = (v_1 + v_2)(w_2 - w_1) \tag{13}$$

Figure 1. Given $A$.
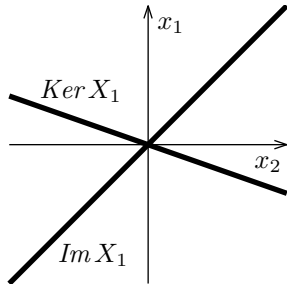


Figure 2. Arbitrary gen. inv. $X$.
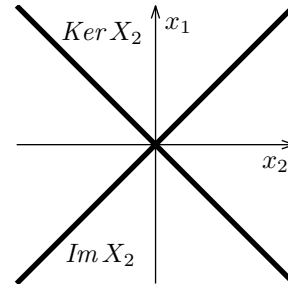


Figure 3. Gen. inv. $X_1$ with arbitrary $Ker\,X_1$.



Figure 4. MP inverse $X_2 = A^\dagger$.

and $S_M \neq 0$ implies $v_1 \neq -v_2$ and $w_1 \neq w_2$ or, in other words, $Im\,X \neq Ker\,A$ and $Ker\,X \neq Im\,A$ for any generalized inverse $X$. The orthogonal projector on $Im\,A$ reads as follows:

$$P_A = \left( \begin{array}{cc} 1/2 & 1/2 \\ 1/2 & 1/2 \end{array} \right).$$

Now

$$X_1 := P_A X = S_M^{-1} \left( \begin{array}{c} (v_1 + v_2)/2 \\ (v_1 + v_2)/2 \end{array} \right) (w_2, -w_1)$$

is the generalized inverse with $Im\,X_1 = Im\,A$ and $Ker\,X_1 = Ker\,X$; see Figure 3. Further

$$X_2 := X_1 P_A = S_M^{-1} \left( \begin{array}{c} (v_1 + v_2)/2 \\ (v_1 + v_2)/2 \end{array} \right) ((w_2 - w_1)/2, (w_2 - w_1)/2)$$

is the generalized inverse with $Im\,X_2 = Im\,A$ and $Ker\,X_2 = Ker\,A$; see Figure 4. Therefore $X_2 = A^\dagger$ by Theorem 2.2 and, moreover due to (13), it follows

$$A^\dagger = \left( \begin{array}{c} 1/2 \\ 1/2 \end{array} \right) (1/2, 1/2) = \left( \begin{array}{cc} 1/4 & 1/4 \\ 1/4 & 1/4 \end{array} \right).$$

The observations above hold in general.

**Theorem 2.3.** *Let $A \in \mathbb{R}^{m \times n}$ be given. Let $X$ be an arbitrary generalized inverse to $A$ and let $P_A$ and $P_{A^\top}$ be the orthogonal projectors on $Im\, A$ and $Im\, A^\top$, respectively. Then*

$$A^\dagger = P_{A^\top} X P_A. \tag{14}$$

*Proof.* Notice that $P_A$, $P_{A^\top}$ are symmetric by Lemma 2.5. We obtain $P_A A = A$ and $A P_{A^\top} = (P_{A^\top} A^\top)^\top = (A^\top)^\top = A$. Now we will verify that $A^\dagger$ defined by (14) satisfies $(11)_1$-$(11)_4$. The first equality $(11)_1$ is straightforward since

$$AA^\dagger A = A P_{A^\top} X P_A A = AXA = A. \tag{15}$$

To prove $(11)_2$ we take arbitrary $x, y \in \mathbb{R}^n$ and consider respective $z_x, z_y \in \mathbb{R}^m$ so that $P_{A^\top} x = A^\top z_x$, $P_{A^\top} y = A^\top z_y$. Then

$$
\begin{aligned}
x^\top A^\dagger A y &= x^\top P_{A^\top} X P_A A y = x^\top P_{A^\top} X A y = z_x^\top A X A y = z_x^\top A y = x^\top P_{A^\top} y \\
&= x^\top A^\top z_y = x^\top A^\top X^\top A^\top z_y = x^\top A^\top X^\top P_{A^\top} y = x^\top (P_{A^\top} X A)^\top y \\
&= x^\top (P_{A^\top} X P_A A)^\top y = x^\top (A^\dagger A)^\top y
\end{aligned}
$$

yields $(11)_2$. The proof of $(11)_3$ is analogous. Finally, let us consider $x \in \mathbb{R}^n$, $y \in \mathbb{R}^m$, and $z_y \in \mathbb{R}^n$ so that $P_A y = A z_y$. We derive

$$
\begin{aligned}
x^\top A^\dagger A A^\dagger y &= x^\top A^\dagger A P_{A^\top} X P_A y = x^\top A^\dagger A P_{A^\top} X A z_y = x^\top A^\dagger A X A z_y \\
&= x^\top A^\dagger A z_y = x^\top A^\dagger P_A y = x^\top P_{A^\top} X P_A P_A y \\
&= x^\top A^\dagger y
\end{aligned}
$$

that proves $(11)_4$.                                                                                    □

**Remark 2.1.** The orthogonal projectors in (14) can be expressed by the SVD (5) as $P_A = I - V_2 V_2^\top$ and $P_{A^\top} = I - U_2 U_2^\top$. It is easily seen that $V_2$ and $U_2$ in $P_A$ and $P_{A^\top}$ may be replaced by arbitrary matrices whose columns form orthogonal bases in $Ker\, A$ and $Ker\, A^\top$, respectively. In the next section, we will assume that the knowledge of such bases is an à-priori information about our problem.

**Remark 2.2.** If $m = n$ and $A$ is symmetric, then (14) simplifies into $A^\dagger = P_A X P_A$. The necessary and sufficient conditions characterizing the MP inverse take the form:

$$A = AA^\dagger A, \quad Im\, A^\dagger = Im\, A, \quad A^\dagger \text{ is symmetric,}$$

or

$$A = AA^\dagger A, \quad Ker\, A^\dagger = Ker\, A, \quad A^\dagger \text{ is symmetric.}$$

**Remark 2.3.** The formula (14) is generalized in Appendix I.

## 3. MP INVERSE IN SADDLE-POINT SYSTEMS

This section deals with solving saddle-point linear systems with singular diagonal blocks by the method combining the Schur complement reduction with orthogonal projectors. Such solution strategy is an algebraic background for different variants of the FETI method [12, 9]. In terms of the standard saddle-point terminology [2] it is the combination of the *range-space method*

and the *null-space method* applied to the primal and the dual saddle-point system, respectively. We shall see that the use of the MP inverse based on (14) is natural and it simplifies both the implementation as well as the analysis.

First of all we introduce notation. Let $\mathbb{V} \subseteq \mathbb{R}^q$ be a subspace. The kernel and the image of any matrix $M \in \mathbb{R}^{p \times q}$ on $\mathbb{V}$ will be denoted by $Ker(M|\mathbb{V})$ and $Im(M|\mathbb{V})$, respectively. If $M$ is symmetric, positive semi-definite (with $p = q$) on $\mathbb{V}$, we will denote the largest eigenvalue on $\mathbb{V}$ by $\sigma_{\max}(M|\mathbb{V})$ and the smallest eigenvalue on $\mathbb{V}$ by $\sigma_{\min}(M|\mathbb{V})$. The spectral condition number of $M$ on $\mathbb{V}$ is defined by

$$\kappa(M|\mathbb{V}) := \frac{\sigma_{\max}(M|\mathbb{V})}{\sigma_{\min}(M|\mathbb{V})}. \tag{16}$$

Moreover, when $\mathbb{V} = \mathbb{R}^q$, we write as before $Ker\,M = Ker(M|\mathbb{V})$, $Im\,M = Im(M|\mathbb{V})$, and $\sigma_{\min}(M) = \sigma_{\min}(M|\mathbb{V})$, $\sigma_{\max}(M) = \sigma_{\max}(M|\mathbb{V})$, $\kappa(M) := \kappa(M|\mathbb{V})$. Let us note that $0 < \sigma_{\min}(M|Im\,M)$, $\sigma_{\max}(M|Im\,M) = \sigma_{\max}(M)$, and $\kappa(M|Im\,M) < +\infty$, if $M$ is the nonzero matrix.

### 3.1. Algorithm

We shall consider the problem to find $(\bar{u}, \bar{\lambda}) \in \mathbb{R}^n \times \mathbb{R}^m$ satisfying:

$$\mathcal{A} \begin{pmatrix} u \\ \lambda \end{pmatrix} = \begin{pmatrix} f \\ g \end{pmatrix} \tag{17}$$

with the saddle-point matrix

$$\mathcal{A} := \begin{pmatrix} A & B^\top \\ B & 0 \end{pmatrix},$$

where $A \in \mathbb{R}^{n \times n}$ is symmetric, positive semi-definite, and singular, $B \in \mathbb{R}^{m \times n}$, $f \in \mathbb{R}^n$, and $g \in \mathbb{R}^m$. We suppose that (17) is uniquely solvable, which is guaranteed by the following necessary and sufficient conditions [17]:

$$Ker\,B^\top = \{0\}, \tag{18}$$
$$Ker\,A \cap Ker\,B = \{0\}. \tag{19}$$

Notice that (18) is the condition on the full row-rank of $B$. Moreover, we assume that an orthonormal basis of $Ker\,A$ is known à-priori and that its vectors are columns of $R \in \mathbb{R}^{n \times l}$, $l = n - r(A)$. Then $P_A = I - RR^\top$ is the orthogonal projector on $Im\,A$ and the MP inverse to $A$ is given by Theorem 2.3 as

$$A^\dagger = (I - RR^\top)X(I - RR^\top), \tag{20}$$

where $X$ is an arbitrary generalized inverse to $A$. Under our assumptions $X$ is easily available by a variant of the Cholesky factorization [6].

The first equation in (17) is satisfied iff

$$f - B^\top \bar{\lambda} \in Im\,A \tag{21}$$

and

$$\bar{u} = A^\dagger(f - B^\top \bar{\lambda}) + R\bar{\alpha} \tag{22}$$

http://hdl.handle.net/10084/94957

for an appropriate $\bar{\alpha} \in \mathbb{R}^l$. To prove this equivalence, it is enough to verify the implication "$\Leftarrow$", since the opposite implication is trivial. Due to (21), there is $v \in \mathbb{R}^n$ such that $Av = f - B^\top \bar{\lambda}$ and, then,

$$A\bar{u} = AA^\dagger(f - B^\top \bar{\lambda}) + AR\bar{\alpha} = AA^\dagger Av = Av = f - B^\top \bar{\lambda}$$

gives the required result. Let us note that $A^\dagger(f - B^\top \bar{\lambda}) \in Im\,A$ and $R\bar{\alpha} \in Ker\,A$. Since $Im\,A$ is the orthogonal complement of $Ker\,A$, $\bar{\alpha}$ is determined uniquely by (22) and, moreover, (21) can be equivalently written as

$$R^\top(f - B^\top \bar{\lambda}) = 0. \tag{23}$$

Further substituting (22) into the second equation in (17) we arrive at

$$-BA^\dagger B^\top \bar{\lambda} + BR\bar{\alpha} = g - BA^\dagger f. \tag{24}$$

Summarizing (24) and (23) we find that the pair $(\bar{\lambda}, \bar{\alpha}) \in \mathbb{R}^m \times \mathbb{R}^l$ satisfies:

$$\mathcal{S} \begin{pmatrix} \lambda \\ \alpha \end{pmatrix} = \begin{pmatrix} d \\ e \end{pmatrix}, \tag{25}$$

where

$$\mathcal{S} := \begin{pmatrix} BA^\dagger B^\top & -BR \\ -R^\top B^\top & 0 \end{pmatrix}$$

is the (negative) *Schur complement* of $A$ in $\mathcal{A}$, $d := BA^\dagger f - g$, and $e := -R^\top f$. As both $\mathcal{S}$ and $\mathcal{A}$ are simultaneously invertible [17], we can compute first $(\bar{\lambda}, \bar{\alpha})$ by solving (25) and then we obtain $\bar{u}$ from (22). Let us note that (25) has formally the same saddle-point structure as that of (17), however, its size is considerably smaller.

Before discussing the solution method for (25) we introduce new notation

$$F := BA^\dagger B^\top, \quad G := -R^\top B^\top$$

which changes (25) into

$$\begin{pmatrix} F & G^\top \\ G & 0 \end{pmatrix} \begin{pmatrix} \lambda \\ \alpha \end{pmatrix} = \begin{pmatrix} d \\ e \end{pmatrix}. \tag{26}$$

Now we shall split (26) using the orthogonal projector $P_G$ on $Ker\,G$. As (19) implies that $G$ is of full row-rank, we can identify $P_G$ with the following matrix:

$$P_G := I - G^\top(GG^\top)^{-1}G.$$

Applying $P_G$ on the first equation in (26) we obtain that $\bar{\lambda}$ satisfies:

$$P_G F\lambda = P_G d, \quad G\lambda = e. \tag{27}$$

In order to arrange (27) as one equation on the vector space $Ker\,G$ we decompose the solution $\bar{\lambda}$ into $\bar{\lambda}_{Im} \in Im\,G^\top$ and $\bar{\lambda}_{Ker} \in Ker\,G$ as

$$\bar{\lambda} = \bar{\lambda}_{Im} + \bar{\lambda}_{Ker}. \tag{28}$$

Since $\bar{\lambda}_{Im}$ is easily available via

$$\bar{\lambda}_{Im} = G^\top(GG^\top)^{-1}e,$$

it remains to show how to get $\bar{\lambda}_{Ker}$. Substituting (28) into (27) we can see that $\bar{\lambda}_{Ker}$ satisfies:

$$P_G F \lambda_{Ker} = P_G(d - F\bar{\lambda}_{Im}), \quad \lambda_{Ker} \in Ker\,G. \tag{29}$$

Let us note that this equation is uniquely solvable, as $P_G F : Ker\,G \mapsto Ker\,G$ is invertible iff $\mathcal{A}$ is invertible [17]. Then, if $\bar{\lambda}$ is known, the solution component $\bar{\alpha}$ is given by

$$\bar{\alpha} = (GG^\top)^{-1}G(d - F\bar{\lambda}). \tag{30}$$

Let us summarize the previous results algorithmically. It turns out to be reasonable to form and store the $l \times m$ matrix $G$ and the $l \times l$ matrix $H := (GG^\top)^{-1}$ because $l$ is usually small (the Cholesky factor of $GG^\top$ may be used instead of $H$). On the other hand, the $m \times m$ matrices $F$ and $P_G$ are not assembled explicitly, since only their matrix-vector products are needed. The actions on $\lambda$ can be evaluated successively as indicated by parentheses on the right hand-sides of

$$F\lambda := B(A^\dagger(B^\top\lambda)) \quad \text{and} \quad P_G\lambda := \lambda - G(H(G^\top\lambda)).$$

In our problems, the actions of $B$ and $B^\top$ are inexpensive due to sparsity of $B$. The actions of $A^\dagger$ are computed by the formula (20) whose implementation is discussed in more details in Section 5.

### ALGORITHMIC SCHEME

Step 1.a:  Assemble $G := -R^\top B^\top$, $H := (GG^\top)^{-1}$, $d := BA^\dagger f - g$, and $e := -R^\top f$.
Step 1.b:  Assemble $\bar{\lambda}_{Im} := G^\top He$.
Step 1.c:  Assemble $\tilde{d} := d - F\bar{\lambda}_{Im}$.
Step 1.d:  Compute $\bar{\lambda}_{Ker}$ by solving $P_G F \lambda_{Ker} = P_G\tilde{d}$ on $Ker\,G$.
Step 1.e:  Assemble $\bar{\lambda} := \bar{\lambda}_{Im} + \bar{\lambda}_{Ker}$.
Step 2:    Assemble $\bar{\alpha} := HG(d - F\bar{\lambda})$.
Step 3:    Assemble $\bar{u} := A^\dagger(f - B^\top\bar{\lambda}) + R\bar{\alpha}$.

### 3.2. Eigenvalue analysis

The key point of the presented algorithm is the equation (29) used in Step 1.d. Its solution can be computed by the projected variant of the conjugate gradient method [12]; see Appendix II. As the rate of convergence of this method is determined by the condition number [14], we shall analyze bounds on the eigenvalues of $P_G F$ on $Ker\,G$.

First note that $P_G F$ is symmetric on $Ker\,G$:

$$\mu^\top P_G F \nu = \mu^\top P_G F P_G \nu = \nu^\top P_G^\top F^\top P_G^\top \mu = \nu^\top P_G F \mu \quad \forall \mu, \nu \in Ker\,G.$$

It is also easy to see that $P_G F$ is positive semi-definite on $Ker\,G$:

$$\mu^\top P_G F \mu = \mu^\top P_G BA^\dagger B^\top P_G \mu \geq 0 \quad \forall \mu \in Ker\,G.$$

As $P_G F$ is invertible on $Ker\,G$, one can deduce that it is positive definite on $Ker\,G$. Bellow we will prove a positive lower bound for the smallest eigenvalue of $P_G F$ on $Ker\,G$. To this end, we will assume that there are constants $0 < c_{A,1} < c_{A,2}$ such that

$$c_{A,1} \leq \sigma_{\min}(A|Im\,A) \quad \text{and} \quad \sigma_{\max}(A) \leq c_{A,2}. \tag{31}$$

Moreover, as the matrix $BB^\top$ is positive definite due to (18), there are constants $0 < c_{B,1} \leq c_{B,2}$ such that

$$c_{B,1} \leq \sigma_{\min}(BB^\top) \quad \text{and} \quad \sigma_{\max}(BB^\top) \leq c_{B,2}. \tag{32}$$

We obtain immediately the following result.

**Lemma 3.7.** *Let $A^\dagger$ be the MP inverse to symmetric, positive semi-definite $A$. Then*

$$c_{A,2}^{-1} \leq \sigma_{\min}(A^\dagger | Im\, A), \quad \sigma_{\max}(A^\dagger) \leq c_{A,1}^{-1}. \tag{33}$$

*Proof:* It follows from the definition (9) of $A^\dagger$, since the nonzero eigenvalues of $A$ are the diagonal entries of $\widehat{\Sigma}$ in the SVD (5), and $Im\, A = Im\, A^\dagger$.                                   □

Now we shall prove the main result of this section.

**Theorem 3.4.** *Let $P_G F$ be the operator of (29). Then*

$$c_{A,2}^{-1} c_{B,1} \leq \sigma_{\min}(P_G F | Ker\, G), \quad \sigma_{\max}(P_G F | Ker\, G) \leq c_{A,1}^{-1} c_{B,2}, \tag{34}$$

*and*

$$\kappa(P_G F | Ker\, G) \leq \frac{c_{B,2}}{c_{B,1}} \cdot \frac{c_{A,2}}{c_{A,1}}. \tag{35}$$

*Proof:* As the proofs of both bounds (34) are analogous, we confine ourself to the lower bound:

$$
\begin{aligned}
\sigma_{\min}(P_G F | Ker\, G) \;=\; & \min_{\substack{\mu \in Ker\, G \\ \mu \neq 0}} \frac{\mu^\top P_G F \mu}{\mu^\top \mu} \;=\; \min_{\substack{R^\top B^\top \mu = 0 \\ \mu \neq 0}} \frac{\mu^\top B A^\dagger B^\top \mu}{\mu^\top \mu} \;=\; \\
=\; & \min_{\substack{R^\top v = 0 \\ v = B^\top \mu \\ \mu \neq 0}} \frac{v^\top A^\dagger v}{v^\top v} \cdot \frac{\mu^\top BB^\top \mu}{\mu^\top \mu} \;\geq\; \\
\geq\; & \min_{\substack{v \in Im\, A \cap Im\, B^\top \\ v \neq 0}} \frac{v^\top A^\dagger v}{v^\top v} \cdot \min_{\substack{\mu \in Ker\, G \\ \mu \neq 0}} \frac{\mu^\top BB^\top \mu}{\mu^\top \mu}.
\end{aligned}
$$

Further using (32),

$$\min_{\substack{\mu \in Ker\, G \\ \mu \neq 0}} \frac{\mu^\top BB^\top \mu}{\mu^\top \mu} = \sigma_{\min}(BB^\top | Ker\, G) \geq \sigma_{\min}(BB^\top) \geq c_{B,1}$$

and, by Lemma 3.7,

$$\min_{\substack{v \in Im\, A \cap Im\, B^\top \\ v \neq 0}} \frac{v^\top A^\dagger v}{v^\top v} \geq \min_{\substack{v \in Im\, A \\ v \neq 0}} \frac{v^\top A^\dagger v}{v^\top v} = \sigma_{\min}(A^\dagger | Im\, A) \geq c_{A,2}^{-1}.$$

Therefore

$$\sigma_{\min}(P_G F | Ker\, G) \geq c_{A,2}^{-1} c_{B,1}$$

that is the lower bound. The inequality (35) follows immediately from (34).            □

**Remark 3.4.** Let the MP inverse $A^\dagger$ in $F$ be replaced by an arbitrary generalized inverse $X$ satisfying solely $AXA = A$. Then Theorem 3.4 remains valid. To prove this result it is enough to show that the eigenvalue bounds (33) are independent on the choice of a generalized inverse:

$$\sigma_{\min}(X|Im\,A) = \min_{\substack{v \in Im\,A \\ v \neq 0}} \frac{v^\top X v}{v^\top v} = \min_{\substack{w = w_0 + w_1 \\ w_0 \in Ker\,A, w_1 \in Im\,A \\ w_1 \neq 0}} \frac{w^\top A X A w}{w^\top A^2 w} = \min_{\substack{w_1 \in Im\,A \\ w_1 \neq 0}} \frac{w_1^\top A w_1}{w_1^\top A^2 w_1}.$$

Therefore

$$\sigma_{\min}(X|Im\,A) = \sigma_{\min}(A^\dagger|Im\,A) \geq c_{A,2}^{-1}$$

and analogously for the largest eigenvalue.

**Remark 3.5.** When $B$ is not full-row rank matrix, then $P_G F$ is singular on $Ker\,G$. In particular, there is $\mu_0 \in Ker\,B^\top$, $\mu_0 \neq 0$, $B^\top \mu_0 = 0$, and $G\mu_0 = R^\top B^\top \mu_0 = 0$, for which

$$\sigma_{\min}(P_G F|Ker\,G) = \frac{\mu_0^\top B A^\dagger B^\top \mu_0}{\mu_0^\top \mu_0} = 0.$$

## 4. APPLICATION IN THE TFETI METHOD

Applying previous results to the saddle-point system (17) arising from the TFETI method [9] we will prove that the condition number of $P_G F$ on $Ker\,G$ does not depend on the size of the problem. First we mention main principles of the TFETI method.

The FETI as well as the TFETI methods belong to the class of non-overlapping domain decomposition methods proposed for the parallel solution of boundary value problems described by elliptic PDEs on a bounded domain $\Omega \subset \mathbb{R}^d$, $d = 1, 2, 3$. Let $L$ denote the diameter of $\Omega$. The idea consists in decomposing $\Omega$ into subdomains $\Omega_k$, $k = 1, \ldots, s$, so that $\overline{\Omega} = \bigcup_{k=1}^s \overline{\Omega}_k$ and $\Omega_k \cap \Omega_l = \emptyset$, $k \neq l$, and considering the PDEs independently in each $\Omega_k$. Therefore the corresponding stiffness matrix $A$ exhibits the block diagonal structure

$$A = diag(A_1, \ldots, A_s) \tag{36}$$

with $A_k \in \mathbb{R}^{N_k \times N_k}$ being symmetric positive semi-definite. Let us note that the number of subdomains $s$ is typically proportional to $(L/H)^d$, where $H$ is the diameter of the largest $\Omega_k$, while the size $N_k$ of $A_k$ corresponds to $(H/h)^d$, where $h$ is the element norm of the finite element approximation. The subdomain interconnectivity is enforced by the second equation in (17) with

$$B = (B_1, \ldots, B_s) \quad \text{and} \quad g = 0, \tag{37}$$

where $B_k \in \mathbb{R}^{m \times N_k}$ and $m$ is proportional to $\sum_{k=1}^s N_k^{(d-1)/d}$. Let us note that finite element nodes shared by more than two subdomains generate usually linearly dependent rows in $B$. In agreement with (18) we will assume that this redundancy is eliminated from $B$ and that the resulting full-row rank matrix is denoted by $B$ again.

The TFETI method enforces also the Dirichlet boundary conditions through the matrix $B$. The main advantage of this strategy is the fact that for each $\Omega_k$ we generate the corresponding block $A_k$ of $A$ as the stiffness matrix to the original PDEs with the pure homogeneous Neumann conditions on the boundary of $\Omega_k$. Consequently, all blocks $A_k$ exhibit the same

kernel dimension, say $\hat{l}$, and their kernel basis may be identified by a mechanical interpretation of the PDEs [9]. In particular, we can assemble the basis for $Ker\,A$ in the matrix $R \in \mathbb{R}^{n \times s\hat{l}}$, $n = \sum_{k=1}^{s} N_k$, with the following block diagonal structure:

$$R = diag(R_1, \ldots, R_s), \tag{38}$$

where $R_k \in \mathbb{R}^{N_k \times \hat{l}}$ may be obtained without any computation (or at negligible cost) as basis vectors of the rigid body motions. As (20) requires orthogonality of $R$, we shall assume that columns of $R_k$ are orthogonalized and that the resulting orthogonal matrix is denoted by $R_k$ again. Let us note that the orthogonalization procedure (if it is necessary) is cheap due to the special structure of $R_k$.

*4.1. Model problem in 1D*

Let $\Omega = (0, L)$, $L > 0$. Let us consider the following problem:

$$-u'' = b \ \text{ in } \Omega, \quad u(0) = 0, \quad u'(L) = 0, \tag{39}$$

where $b \in C(\Omega)$. Let all subdomains $\Omega_k$ of $\Omega$ be of the same lengths $H = L/s$ so that $\Omega_k = ((k-1)H, kH)$, $k = 1, \ldots, s$. On each $\Omega_k$ we consider an equidistant partition with $N$ nodes so that $h = H/(N-1)$. Note that the decomposition parameter $H$ and the discretization parameter $h$ satisfy:

$$N = 1 + \frac{H}{h}. \tag{40}$$

The approximation of (39) based on the TFETI method with the linear finite elements leads to the blocks in (36) given by $A_k = A(h, N) \in \mathbb{R}^{N \times N}$, where

$$A(h, N) = \frac{1}{h} \begin{pmatrix} 1 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 1 \end{pmatrix}. \tag{41}$$

In each block $B_k \in \mathbb{R}^{m \times N}$ of (37) there are at most two nonzero entries, i.e. "1" at the first position of the $k$th row and "$-1$" at the last position of the $(k+1)$th row (note that $B_s$ contains only one "1" at the beginning of the last row due to the imposed Neumann boundary condition). Finally, all blocks $R_k \in \mathbb{R}^{N \times 1}$ of (38) read as follows:

$$R_k = \frac{1}{\sqrt{N}}(1, \ldots, 1)^\top. \tag{42}$$

**Lemma 4.8.** *The eigenvalues $\sigma_j = \sigma_j(h, N)$ of $A(h, N)$ read as follows:*

$$\sigma_j = \frac{1}{h}(2 - 2\cos\theta_j), \quad \theta_j = \frac{j\pi}{N}, \quad j = 0, 1, \ldots, N-1.$$

*Proof:* Using the standard trigonometric formulas one can verify that $v_j = (\cos(i + \frac{1}{2})\theta_j)_{i=0}^{N-1}$ is the eigenvector corresponding to $\sigma_j$; see [24]. $\qquad\square$

**Theorem 4.5.** *Let $N \geq 4$. Let $P_G F$ be the operator of* (29) *given by* (36)-(38) *arising from the TFETI method applied to* (39). *It holds:*

$$\frac{h}{4} \leq \sigma_{\min}(P_G F | Ker\, G),$$

$$\sigma_{\max}(P_G F | Ker\, G) \leq \frac{24}{11\pi^2} \cdot \frac{(h+H)^2}{h},$$

$$\kappa(P_G F | Ker\, G) \leq \frac{96}{11\pi^2} \left(1 + \frac{H}{h}\right)^2.$$

*Proof:* We estimate the constants from (31) and (32). As $BB^\top = diag(1, 2, \ldots, 2) \in \mathbb{R}^{s \times s}$, we obtain immediately $c_{B,1} = 1$ and $c_{B,2} = 2$. As all diagonal blocks of $A$ are $A(h, N)$, it follows from Lemma 4.8 that

$$\sigma_{\min}(A | Im\, A) = \sigma_1 \quad \text{and} \quad \sigma_{\max}(A) = \sigma_{N-1}.$$

Therefore we can take $c_{A,2} = 4/h$. To estimate $c_{A,1}$ we use the Taylor expansion for $\cos \theta_1$:

$$\sigma_1 = \frac{2}{h}(\theta_1^2/2! - \theta_1^4/4! + \theta_1^6/6! - \theta_1^8/8! - \ldots).$$

As $N \geq 4$, we have $\theta_1 \leq \pi/4 < 1$ so that $\theta_1^{2i}/(2i)! - \theta_1^{2(i+1)}/(2(i+1))! \geq 0$ and therefore

$$\sigma_1 \geq \frac{2}{h}(\theta_1^2/2! - \theta_1^4/4!).$$

Further $\theta_1^2 \geq \theta_1^4$ implies

$$\sigma_1 \geq \frac{11}{12h}\theta_1^2 = c_{A,1}.$$

Substituting $\theta_1 = \pi/N$ and (40) we arrive at

$$c_{A,1} = \frac{11\pi^2 h}{12(h+H)^2}.$$

The rest of the proof consists in using Theorem 3.4. $\qquad\qquad\qquad\qquad\qquad\square$

**Remark 4.6.** For $N = 3$, we can take $c_{A,1} = \sigma_1 = 1/h$.

*4.2. Model problem in 2D*

Let $\Omega = (0, L_x) \times (0, L_y)$, $L_x, L_y > 0$. Let us consider the Poisson problem:

$$-\Delta u = b \ \text{ in } \Omega, \quad u = 0 \ \text{ on } \gamma_d, \quad \frac{\partial u}{\partial n} = 0 \ \text{ on } \gamma_n, \qquad (43)$$

where $\gamma_d = \{0\} \times (0, L_y)$, $\gamma_n = \partial\Omega \setminus \overline{\gamma}_d$, $b \in C(\Omega)$, and $n$ in $\frac{\partial u}{\partial n}$ denotes the unit outer normal vector to the boundary $\partial\Omega$. Let all subdomains of $\Omega$ be rectangles $\Omega_k = \Omega_{k_x} \times \Omega_{k_y}$, where $\Omega_{k_z} = ((k_z - 1)H_z, k_z H_z)$, $k_z = 1, \ldots, s_z$, $H_z = L_z/s_z$ for $z = x, y$. The number of $\Omega_k$ is $s = s_x s_y$ and the correspondence between $k_x$, $k_y$ and $k$ is given by $k = k_x + (k_y - 1)s_x$. Let us

construct equidistant partitions of the sides of $\Omega_k$ with the same stepsizes $h_x = H_x/(N_x - 1)$ and $h_y = H_y/(N_y - 1)$ for all $k$. Thus, each $\Omega_k$ is partitioned by $N = N_x N_y$ nodes into $(N_x - 1)(N_y - 1)$ rectangles. Finally, we assume that every rectangle is cut by its diagonal into two triangles. On this triangulation of $\Omega_k$ we define the finite-element space of continuous piecewise linear functions that is used to approximate the solution to (43) by the TFETI method. The blocks $A_k$ in (36) are given by

$$A_k = A_x \otimes D_y + D_x \otimes A_y, \tag{44}$$

where $A_z = A(h_z, N_z) \in \mathbb{R}^{N_z \times N_z}$ is defined by (41), $D_z = h_z I_z \in \mathbb{R}^{N_z \times N_z}$ is diagonal, $z = x, y$, and $\otimes$ stands for the Kronecker tensor product of matrices. The nonzero entries of blocks $B_k$ in (37) are "1" and "$-1$" at appropriate positions corresponding to the nodes lying on the boundaries $\partial \Omega_k$ (the signs reflect an orientation of the outer normal vector). Recall that we assume full-row rank $B$ without redundant rows. In order to simplify the next presentation we assume that, in addition, the rows of $B$ are orthogonal (due to an orthogonalization procedure at negligible cost). Finally, note that the blocks $R_k \in \mathbb{R}^{N \times 1}$ in (38) are given by (42) again.

**Theorem 4.6.** *Let $N_x, N_y \geq 4$ and denote $\delta = h_x/h_y$. Let $P_G F$ be the operator of* (29) *given by* (36)-(38) *arising from the TFETI method applied to* (43). *It holds:*

$$\frac{1}{4} \left( \delta^{-1} + \delta \right)^{-1} \quad \leq \quad \sigma_{\min}(P_G F | Ker\, G),$$

$$\sigma_{\max}(P_G F | Ker\, G) \quad \leq \quad \frac{12}{11\pi^2} \cdot \max \left\{ \delta \left( 1 + \frac{H_x}{h_x} \right)^2, \delta^{-1} \left( 1 + \frac{H_y}{h_y} \right)^2 \right\},$$

$$\kappa(P_G F | Ker\, G) \quad \leq \quad \frac{48}{11\pi^2} \cdot \max \left\{ (1 + \delta^2) \left( 1 + \frac{H_x}{h_x} \right)^2, (1 + \delta^{-2}) \left( 1 + \frac{H_y}{h_y} \right)^2 \right\}.$$

*Proof:* The proof is analogous as for Theorem 4.5. Now $c_{B,1} = c_{B,2} = 1$, as $B$ is orthogonal. It is well-known from the Kronecker product structure of (44) that each eigenvalue $\sigma(A_k)$ of $A_k$ is of the form $\sigma(A_k) = \sigma(A_x)\sigma(D_y) + \sigma(D_x)\sigma(A_y)$, where $\sigma(A_z)$ and $\sigma(D_z)$ are the eigenvalues of $A_z$ and $D_z$, $z = x, y$, respectively [14]. Since all eigenvalues of $D_z$ are $h_z$, $z = x, y$, and the eigenvalues of $A_z$ are given by Lemma 4.8 as $\sigma_{j,z} = \sigma_j(h_z, N_z)$, $j = 0, 1, \ldots, N_z - 1$, $z = x, y$, we obtain

$$\sigma_{\min}(A | Im\, A) = \min\{h_y \sigma_{1,x}, h_x \sigma_{1,y}\}, \quad \sigma_{\max}(A) = h_y \sigma_{N_x - 1, x} + h_x \sigma_{N_y - 1, y}.$$

Applying the same bounds as in the proof of Theorem 4.5, we get

$$\sigma_{\max}(A) \quad \leq \quad h_y \frac{4}{h_x} + h_x \frac{4}{h_y} = 4(\delta^{-1} + \delta) \quad =: \quad c_{A,2},$$

$$\sigma_{\min}(A | Im\, A) \quad \geq \quad \frac{11\pi^2}{12} \cdot \min \left\{ h_y \frac{h_x}{(h_x + H_x)^2}, h_x \frac{h_y}{(h_y + H_y)^2} \right\} \quad =$$

$$= \quad \frac{11\pi^2}{12} \cdot \min \left\{ \delta^{-1} \left( 1 + \frac{H_x}{h_x} \right)^{-2}, \delta \left( 1 + \frac{H_y}{h_y} \right)^{-2} \right\} \quad =: \quad c_{A,1}.$$

The rest consists in using Theorem 3.4.          $\square$

**Remark 4.7.** Denote $h = (h_x^2 + h_y^2)^{1/2}$ and $H = (H_x^2 + H_y^2)^{1/2}$. If $h_x = h_y$ and $H_x = H_y$, then $\delta = 1$ and the bound on $\kappa(P_G F | Ker\, G)$ is the same as for the model problem in 1D.

**Remark 4.8.** The results of Theorem 4.5 and Theorem 4.6 can be improved by the analysis of [3, 4, 5, 12].

## 5. NUMERICAL EXPERIMENTS

Numerical experiments will illustrate the above theoretical results for more complex problems arising from linear elasticity in 3D. We will solve (29) by the projected conjugate gradient method [12] (ProjCGM) with the relative terminating precision $1e$-4 for which we will observe the number of iterations (**iter**). A short description of the ProjCGM algorithm is presented in Appendix II. We will experimentally asses sensitivities of computations with respect to the choice of generalized inverses and spectral properties of off-diagonal blocks. We use our parallel MATLAB library MatSol [18].



Figure 5. Geometry of the model problem.

The model problem is the steel cubic body $\Omega \subset \mathbb{R}^3$ as depicted in Figure 5 with the edge length $a = 10$[mm] and the curved top face with the radius $r = 10^4$[mm]. Elastic properties of $\Omega$ are described by the Lamè PDEs with the Young modulus $E = 2e5$[MPa] and the Poisson ratio $\nu = 0.35$. The body is fixed in all directions along the left face and loaded by the vertical traction $p = -2000$[MPa] along the curved top face. The finite element discretization uses uniform trilinear bricks with lexicographic ordering of nodes so that the last ($M$-th) node is on the curved top edge as it is seen in Figure 5. Each subdomain $\Omega_k$ in the TFETI domain decomposition of $\Omega$ exhibits six rigid body modes (three translations and three rotations) so that $\dim Ker\, A_k = 6$ for each block $A_k$ in (36). The blocks $R_k$ in (38) may be assembled by the coordinates of the finite element nodes of $\overline{\Omega}_k$ [9].

*5.1. Actions of generalized inverses*

Before giving numerical experiments we describe generalized inverses $X_k$ to $A_k$ that we use in our computations. Note that the implementation of the algorithm does not require to assemble $X_k$ explicitly. What is only needed it is the action of $X_k$ on a vector.

For the sake od simplicity we omit the index $k$ so that we consider a symmetric, positive semi-definite $A \in \mathbb{R}^{N \times N}$ with $\dim Ker A = 6$. First of all the preprocessing computation identifies the nonsingular part of $A$ by using the LU-factorization $A = LU$. As A is symmetric, we can write $A = LDL^\top$, where the diagonal matrix $D$ is given by the pivots of the LU-factorization, i.e. $D = diag(U)$. In $D$ we find six vanishing (critically small) diagonal entries and introduce the permutation matrix $P_D$ so that

$$P_D D P_D^\top = \left( \begin{array}{cc} D_{11} & 0 \\ 0 & 0 \end{array} \right),$$

where $D_{11} \in \mathbb{R}^{(N-6) \times (N-6)}$ is the nonsingular part of $D$. The analogous permutation of $A$ detects the nonsingular part $A_{11} \in \mathbb{R}^{(N-6) \times (N-6)}$ of $A$ via

$$P_D A P_D^\top = \left( \begin{array}{cc} A_{11} & A_{12} \\ A_{21} & A_{22} \end{array} \right). \tag{45}$$

Now we can define the generalized inverse $X$ to $A$ by

$$X = P_D^\top \left( \begin{array}{cc} A_{11}^{-1} & 0 \\ 0 & 0 \end{array} \right) P_D$$

ant its action may be computed by

$$X = P_D^\top \left( \begin{array}{cc} U_{11}^{-1} L_{11}^{-1} & 0 \\ 0 & 0 \end{array} \right) P_D, \tag{46}$$

where $A_{11} = L_{11} U_{11}$ is the LU-factorization of $A_{11}$. Recall that the action of the MP-inverse $A^\dagger$ may be obtained using (46) in (20), where $R$ corresponds to the basis vectors of the rigid body motions.

The zero pivots of $D$ identify the so-called *fixing DOFs* in the finite element nodes that prevent the (floating) body from rigid body motions. The typical configuration of fixing DOFs for $\Omega$ with the planar and curved top face (i.e. with $r = \infty$ and $r < \infty$) is depicted in Figure 6 and Figure 7; we call them *Configuration 1* and *Configuration 2*, respectively. Although the geometry of our model problem uses the curved top face, its radius is so big ($r = 10^4$[mm]) that the positions of fixing DOFs coincide naturally with Configuration 1. Let us note that the correct determination of fixing DOFs is a sensitive problem especially for bodies with a complicated geometry or with a composite material structure. In order to simulate an effect of this sensitivity, we will try to determine for our geometry a generalized inverse, say $\widetilde{X}$, via Configuration 2 (i.e., we set $P_D$ in (45) using the positions of the fixing DOFs in Configuration 2 and, then, we define $\widetilde{X}$ by the right-hand side of (46)). Note that the fixing DOFs in Configuration 2 do not determine any generalized inverse for the planar top face ($r = \infty$), since the rotation of the body $\Omega$ is allowed and, consequently, $A_{11}$ in (45) must be singular.
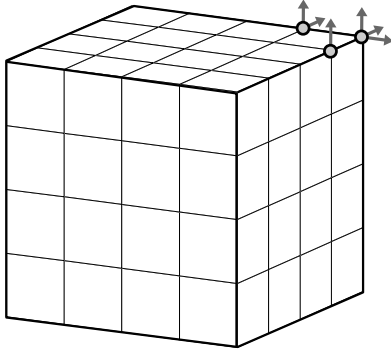
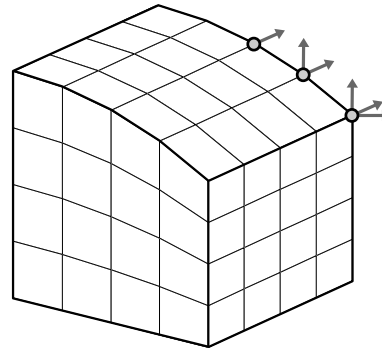Figure 6. Configuration 1 of fixing DOFs.



Figure 7. Configuration 2 of fixing DOFs.

### 5.2. Correct and incorrect generalized inverses

All computations are performed with fixed $H/h = 10$ for various numbers of subdomains $s$, where $H$ and $h$ stand for the decomposition and discretization parameters, respectively, as in Section 4. The symbols $n$, $m$, and $l$ denote the number of primal unknowns, dual unknowns, and the kernel-space dimension, respectively, as in Section 3.

The first numerical experiment is carried out using the correct generalized inverse $X$, i.e. using Configuration 1. The row of Table I labeled by $X$ summarizes the characteristics of solving (29) with $A^\dagger$ replaced by $X$. Comparing with the next row computed by $A^\dagger$, we can conclude in agreement with predictions of the theory that the number of iterations is independent on both the size of the saddle-point system as well as the choice of the generalized inverse.

Table I. Iterations and solution times; **iter**(time in seconds).

| $s$ | 1 | 27 | 125 | 343 | 729 |
|---|---|---|---|---|---|
| $n$ | 3,993 | 107,811 | 499,125 | 1,369,599 | 2,910,897 |
| $m$ | 363 | 21,321 | 108,975 | 310,989 | 675,027 |
| $l$ | 6 | 162 | 750 | 2,058 | 4,374 |
| $X$ | **11**(0.79) | **17**(4.05) | **17**(23.98) | **17**(114.74) | **17**(551.51) |
| $A^\dagger$ | **11**(0.83) | **17**(3.28) | **17**(23.50) | **17**(111.38) | **17**(623.59) |
| $\widetilde{X}$ | **12**(0.89) | **32**(7.53) | **29**(35.27) | **27**(126.19) | **27**(607.27) |
| $P_A\widetilde{X}P_A$ | **11**(0.84) | **25**(5.34) | **26**(31.44) | **25**(133.20) | **25**(584.08) |

The second numerical experiment tests $\widetilde{X}$ determined by Configuration 2. The characteristics of solving (29) obtained by $\widetilde{X}$ and $P_A\widetilde{X}P_A$ with $P_A = I - RR^\top$ are reported in the last two rows of Table I. Although the ProjCGM iterations are reasonable terminated in all cases, the solutions computed by $\widetilde{X}$ are not sufficiently accurate, as it is seen from the total displacements depicted in Figure 8, 9 (scaled 4×) and from the constraint errors summarized in Table II. The explanation is simple: the matrix $\tilde{X}$ is not any generalized inverse. If it would be, then $A^\dagger$

and $P_A\tilde{X}P_A$ would coincide and the numbers of iterations would be the same. Among others, this numerical experiment illustrates the stabilization effect of the projection formula (20): although $\tilde{X}$ is not any generalized inverse, $P_A\tilde{X}P_A$ is an approximation of $A^\dagger$ (but not exactly $A^\dagger$) that leads to the sufficiently accurate results. Moreover, the condition number $\kappa(P_A\tilde{X}P_A)$ is less than $\kappa(\tilde{X})$.

Table II. Constraint errors; $\|Bu - g\|/\|u\|$ with $g = 0$.

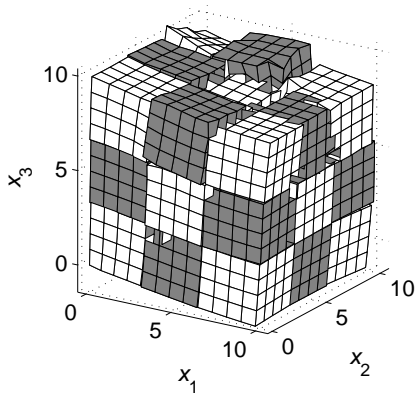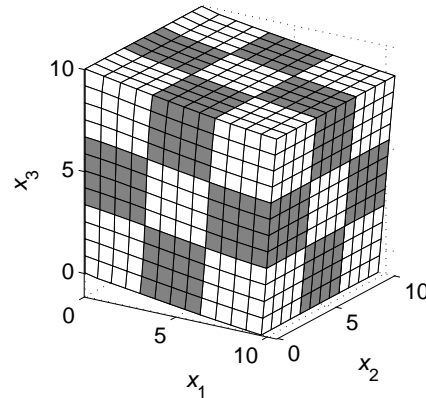| $n$ | $X$ | $A^\dagger$ | $\widetilde{X}$ | $P_A\widetilde{X}P_A$ |
|---|---|---|---|---|
| 3,993 | 4.400e-06 | 4.400e-06 | 0.034 | 1.646e-05 |
| 107,811 | 3.412e-05 | 3.413e-05 | 2.519 | 4.936e-05 |
| 499,125 | 4.788e-05 | 4.788e-05 | 2.261 | 3.613e-05 |
| 1,369,599 | 4.933e-05 | 4.933e-05 | 7.282 | 6.775e-05 |
| 2,910,897 | 5.311e-05 | 5.311e-05 | 55.921 | 7.129e-05 |



Figure 8. Total displacements for $\widetilde{X}$.



Figure 9. Total displacements for $P_A\widetilde{X}P_A$.

### 5.3. Effect of the off-diagonal block

In this example we explore an influence of the off-diagonal block $B$ of the saddle-point system (17) on the behavior of ProjCGM iterations. If the bounds are tight in (32), then Theorem 3.4 gives

$$\kappa(P_G F|Ker G) \leq \frac{c_{A,2}}{c_{A,1}} \cdot \kappa(BB^\top).$$

Since the matrix $B$ glues the subdomain solutions of PDE problems, the spectral properties of $BB^\top$ obviously depend on a decomposition geometry. In order to demonstrate this fact, we denote by $n_x$, $n_y$, and $n_z$ the number of segments in the $x$, $y$, and $z$ directions, respectively,

*Numer. Linear Algebra Appl.* 2011; **00**:1–20

that we use to divide the cube $\Omega$ into the (box) subdomains $\Omega_k$. The matrix $BB^\top$ is (after an appropriate permutation of rows in $B$) block-diagonal. The size of its blocks is given by multiplicities of nodes shared by more subdomains $\Omega_k$. In our geometry one node may belong up to eight subdomains so that the blocks of $BB^\top$ are (after eliminating redundancy from $B$) at most of the seventh order. Table III shows that $BB^\top$ exhibits at most seven different eigenvalues and, in advance, the condition number $\kappa(BB^\top)$ is independent on $H$ and $h$. In Table IV we present the numbers of iterations **iter** for various decompositions. We denote by ORTH(+/-) computations with/without the orthogonalization. Analogously, PREC(+/-) denotes computations with/without preconditioning of the ProjCGM iterations. We use the well-known lumped type preconditioner $\overline{F^{-1}} = BAB^\top$ to $F$ [12, 13] that approximates the inverse to $F$ via the MP-inverse to $A^\dagger$, i.e., via $A$. The column (a) of the table shows that the numbers of iterations are obviously in agreement with the value of $\kappa(BB^\top)$, when the preconditioner and the orthogonalization are not used. The columns (b) and (d) indicate that the orthogonalization of $B$ is necessary for a favorable effect of the preconditioner. The reason is the fact that $B^\top$ in $\overline{F^{-1}}$ plays, in this case, the role of the MP-inverse $B^\dagger$. Finally, the column (c) indicates that the number of iterations may not depend on the number of subdomains (for larger problems), if $B$ is orthogonalized.

Table III. Condition number and eigenvalues of $BB^\top$.

| $n_x \times n_y \times n_z$ | $\kappa(BB^\top)$ | eigenvalues |
|:---:|:---:|:---:|
| $1 \times 1 \times 1$ | 1 | $\{1\}$ |
| $2 \times 1 \times 1$ | 1 | $\{1\}$ |
| $1 \times 2 \times 1$ | 5.8284 | $\{0.2929, 1, 1.7071\}$ |
| $1 \times 1 \times 2$ | 5.8284 | $\{0.2929, 1, 1.7071\}$ |
| $2 \times 2 \times 1$ | 5.8284 | $\{0.2929, 1, 1.7071\}$ |
| $2 \times 1 \times 2$ | 5.8284 | $\{0.2929, 1, 1.7071\}$ |
| $1 \times 2 \times 2$ | 25.274 | $\{0.07612, 0.2929, 0.6173, 1, 1.3827, 1.7071, 1.9239\}$ |
| $2 \times 2 \times 2$ | 25.274 | $\{0.07612, 0.2929, 0.6173, 1, 1.3827, 1.7071, 1.9239\}$ |
| $k \times k \times k,\ k > 2$ | 25.274 | $\{0.07612, 0.2929, 0.6173, 1, 1.3827, 1.7071, 1.9239\}$ |

## 6. CONCLUSIONS AND COMMENTS

We have analyzed the algorithm for solving saddle-point linear systems with singular diagonal blocks that combines the Schur complement reduction with the null-space method. The resulting dual equation is solved by the projected conjugate gradient algorithm. This solution strategy is the background for the classical FETI domain decomposition method [12] and its variants. Since the number of ProjCGM iterations depends on conditioning of the problem, we have derived the bound on the condition number of the corresponding dual operator. Using this result we have discussed the role of the choice of generalized inverses and the effect of conditioning of the off-diagonal block of the saddle-point system. Moreover, we have proved for simple model elliptic boundary value problems that the number of iterations required for the solution with a given accuracy in the TFETI variant [9] of the FETI method does not depend on the discretization and decomposition parameters. The numerical examples confirmed the

Table IV. ProjCGM iterations for $H/h = 5$ and $r = \infty$.

| $n_x \times n_y \times n_z$ | $\kappa(BB^\top)$ | (a) ORTH (-) PREC (-) | (b) ORTH (-) PREC (+) | (c) ORTH (+) PREC (-) | (d) ORTH (+) PREC (+) |
|---|---|---|---|---|---|
| $1 \times 1 \times 1$ | 1 | **15** | **8** | **15** | **8** |
| $2 \times 1 \times 1$ | 1 | **16** | **13** | **16** | **13** |
| $1 \times 2 \times 1$ | 5.8284 | **32** | **35** | **27** | **14** |
| $1 \times 1 \times 2$ | 5.8284 | **30** | **32** | **27** | **11** |
| $2 \times 2 \times 1$ | 5.8284 | **32** | **40** | **27** | **16** |
| $2 \times 1 \times 2$ | 5.8284 | **33** | **36** | **28** | **14** |
| $1 \times 2 \times 2$ | 25.274 | **39** | **91** | **29** | **15** |
| $2 \times 2 \times 2$ | 25.274 | **35** | **77** | **25** | **11** |
| $3 \times 3 \times 3$ | 25.274 | **38** | **92** | **27** | **12** |
| $4 \times 4 \times 4$ | 25.274 | **40** | **111** | **28** | **11** |

theoretical results experimentally for more complicated linear elasticity problems.

Although the theoretical analysis prove that the choice of a generalized inverse does not influence the ProjCGM iterations, the numerical experiments showed that the formula (14) defining the MP inverse may stabilize computations.

## APPENDIX I

We will generalize the formula (14) in the sense that we will show how to obtain the generalized inverse for which the kernel-space and the image-space are prescribed. For the sake of simplicity we confine ourself to symmetric matrices.

Let $A \in \mathbb{R}^{n \times n}$ be symmetric and let $R \in \mathbb{R}^{n \times l}$ be a matrix whose columns are the orthonormal basis of $Ker\,A$. Let us consider the following saddle-point problem: find the pair $(\bar{u}, \bar{\lambda}) \in \mathbb{R}^n \times \mathbb{R}^l$ solving

$$\left( \begin{array}{cc} A & M \\ N^\top & 0 \end{array} \right) \left( \begin{array}{c} u \\ \lambda \end{array} \right) = \left( \begin{array}{c} f \\ 0 \end{array} \right) \tag{47}$$

with $M, N \in \mathbb{R}^{n \times l}$ so that the respective products $R^\top M, R^\top N \in \mathbb{R}^{l \times l}$ are nonsingular and $f \in \mathbb{R}^n$. The first equation in (47) implies $\bar{\lambda} = (R^\top M)^{-1} R^\top f$ and then $A\bar{u} = (I - M(R^\top M)^{-1} R^\top)f$. Therefore

$$\bar{u} = X(I - M(R^\top M)^{-1} R^\top)f + R\bar{\alpha}, \tag{48}$$

where $X \in \mathbb{R}^{n \times n}$ is an arbitrary generalized inverse to $A$ and $\bar{\alpha} \in \mathbb{R}^l$. Substituting (48) in the second equation in (47) we arrive at

$$\bar{\alpha} = -(N^\top R)^{-1} N^\top X(I - M(R^\top M)^{-1} R^\top)f. \tag{49}$$

Simple manipulations with (49) and (48) give

$$\bar{u} = Yf, \tag{50}$$

where

$$Y = P_N^\top X P_M \tag{51}$$

and

$$P_M = I - M(R^\top M)^{-1} R^\top, \quad P_N = I - N(R^\top N)^{-1} R^\top.$$

It is easily seen that $P_M$, $P_N$ are the projectors (not necessarily orthogonal) so that $Im\, P_M = Im\, P_N = Im\, A$ and $Ker\, P_M = Im\, M$, $Ker\, P_N = Im\, N$, respectively. As $P_M A = A$ and $A P_N^\top = A$, we obtain

$$AYA = A P_N^\top X P_M A = AXA = A.$$

Thus, $Y$ is the generalized inverse to $A$ with $Ker\, Y = Im\, M$ and $Im\, Y = Ker\, N^\top$.

Our final remark deals with a mechanical interpretation of the MP inverse. Let $A$ in (47) be the stiffness matrix arising from the finite element approximation of a linearly elastic body $\Omega \subset \mathbb{R}^3$ with the pure Neumann boundary condition (zero surface traction) as introduced in Section 5. Let us choose $M = N = R$ in (47). Then $Y$ in (51) is the MP-inverse to $A$ and the first component of the solution to (47) is given by

$$\bar{u} = A^\dagger f.$$

In other words, since $R^\top \bar{u} = 0$, the MP-inverse $A^\dagger$ to $A$ determines the displacement of the floating body $\Omega$ so that its rigid body modes are prohibited (in average).

## APPENDIX II

Let us introduce the projected conjugate gradient method with preconditioning (ProjCGM) [12] that we use for computing $\bar{\lambda}_{Ker}$ in Step 1.d of Algorithmic scheme. Thus we want to compute $\bar{\lambda}_{Ker}$ by solving the system $P_G F \lambda_{Ker} = P_G \tilde{d}$ on $Ker\, G$ with the lumped type preconditioner $\overline{F^{-1}}$ to $F$.

<u>ALGORITHM PROJCGM</u>

1. Initialize

$$r^0 = \tilde{d}, \ \ \lambda_{Ker}^0 = 0.$$

2. Iterate $k = 1, 2, ...,$ until convergence

   $Project\ w^{k-1} = P_G r^{k-1}.$

   $Precondition\ z^{k-1} = \overline{F^{-1}} w^{k-1}.$

   $Project\ y^{k-1} = P_G z^{k-1}.$

   $\beta^k = (y^{k-1})^\top w^{k-1} / (y^{k-2})^\top w^{k-2}; \qquad (\beta^1 = 0).$

   $p^k = y^{k-1} + \beta^k p^{k-1}; \qquad\qquad (p^1 = y^0).$

   $\alpha^k = (y^{k-1})^\top w^{k-1} / (p^k)^\top F p_k.$

   $\lambda_{Ker}^k = \lambda_{Ker}^{k-1} + \alpha^k p^k.$

   $r^k = r^{k-1} - \alpha^k F p^k.$

3. $\bar{\lambda}_{Ker} = \lambda_{Ker}^k.$

The generalization for non-symmetric systems may be found in [17, 21].

## *Acknowledgement*

## REFERENCES

1. Ben-Israel A, Greville T. *Generalized inverses: theory and applications*, (2nd edn). Springer: New York, 2003.
2. Benzi M, Golub GH, Liesen J. Numerical solution of saddle point systems. *Acta Numerica* 2005; **14**:1–137.
3. Bramble JH, Pasciak JE, Schatz A. The construction of preconditioners for elliptic problems by substructuring I. *Mathematics of Computation* 1986; **47**(175):103–134.
4. Bramble JH, Pasciak JE, Schatz A. The construction of preconditioners for elliptic problems by substructuring II. *Mathematics of Computation* 1987; **49**(179):1–16.
5. Brenner SC. The condition number of the Schur complement in domain decomposition. *Numerische Mathematik* 1999; **83**:187–203.
6. Brzobohatý T, Dostál Z, Kozubek T, Markopoulos A, Kovář P. Cholesky–SVD decomposition with fixing nodes to stable computation of a generalized inverse of the stiffness matrix of a floating structure. Submitted to *International Journal for Numerical Methods in Engineering* 2010.
7. Campbell SL, Meyer CD. *Generalized inverses of linear transformations*. SIAM: Philadelphia, 2009.
8. Dostál Z. *Optimal quadratic programming algorithms: with applications to variational inequalities*, Springer: New York, 2009.
9. Dostál Z, Horák D, Kučera R. Total FETI - an easier implementable variant of the FETI method for numerical solution of elliptic PDE. *Communications in Numerical Methods in Engineering* 2006; **22**(12):1155–1162.
10. Farhat C, Gèrardin M. On the general solution by a direct method of a large scale singular system of linear equations: application to the analysis of floating structures. *International Journal for Numerical Methods in Engineering* 1998; **4**(41):675-696.
11. Farhat C, Lesoinne M, LeTallec P, Pierson K, Rixen D. FETI-DP: a dual-primal unified FETI method. I. A faster alternative to the two-level FETI method. *International Journal for Numerical Methods in Engineering* 2001; **50**:1523–1544.
12. Farhat C, Mandel J, Roux FX. Optimal convergence properties of the FETI domain decomposition method. *Computer Methods in Applied Mechanics and Engineering* 1994; **115**:367–388.
13. Fragakis Y. A study on the lumped preconditioner and memory requirements of FETI and related primal domain decomposition methods. *International Journal for Numerical Methods in Engineering* 2008; **13**(73):1865–1884.
14. Golub GH, Van Loan CF. *Matrix computations* (3th edn). The Johns Hopkins University Press: Baltimore, 1996.
15. Gan G. On the relation between Moore's and Penrose's conditions. *International Journal of Mathematics and Mathematical Sciences* 2002; **30**(8):505-509.
16. Haslinger J, Kozubek T, Kučera R. Fictitious domain formulation of unilateral problems: analysis and algorithms. *Computing* 2009; **84**(1-2):69–96.
17. Haslinger J, Kozubek T, Kučera R, Peichl G. Projected Schur complement method for solving non-symmetric saddle-point systems arising from fictitious domain approach. *Numerical Linear Algebra with Applications* 2007; **14**(9):713–739.
18. Kozubek T, Markopoulos A, Brzobohatý T, Kučera R, Vondrák V, Dostál Z. MatSol - MATLAB efficient solvers for problems in engineering. *http://matsol.vsb.cz*
19. Klawonn A, Widlund OB, Dryja M. Dual-primal FETI methods for three-dimensional elliptic problems with heterogeneous coefficients. *SIAM Journal on Numerical Analysis* 2002; **40**:159–179.
20. Moore EH. On the reciprocal of the general algebraic matrix. *Bulletin of the American Mathematical Society* 1920; **21**:394-395.
21. Orban D. *Projected Krylov methods for unsymmetric augmented systems*. Cahiers du GERARD G-2008-46, GERARD, Montreal, Canada, 2008.
22. Papadrakakis M, Fragakis Y. An integrated geometric–algebraic method for solving semi-definite problems in structural mechanics. *Computer Methods in Applied Mechanics and Engineering* 2001; **190**:6513–6532.

23. Penrose R. A generalized inverse for matrices. *Mathematical Proceedings of the Cambridge Philosophical Society* 1955; **51**:406-413.
24. Strang G. The discrete cosine transform. *SIAM Review* 1999; **41**:135-147.
25. Shinozaki N, Sibuya M, Tanabe K. Numerical algorithms for the Moore-Penrose inverse of a matrix: Direct methods. *Annals of the Institute of Statistical Mathematics* 1972; **24**(1):193–203.