

# ADAPTIVE MULTILEVEL METHODS FOR OBSTACLE PROBLEMS\*

R. H. W. HOPPE† AND R. KORNUBER‡

*Dedicated to Professor R. Bulirsch on the occasion of his 60th birthday.*

**Abstract.** The authors consider the discretization of obstacle problems for second-order elliptic differential operators by piecewise linear finite elements. Assuming that the discrete problems are reduced to a sequence of linear problems by suitable active set strategies, the linear problems are solved iteratively by preconditioned conjugate gradient iterations. The proposed preconditioners are treated theoretically as abstract additive Schwarz methods and are implemented as truncated hierarchical basis preconditioners. To allow for local mesh refinement semilocal and local a posteriori error estimates are derived, providing lower and upper estimates for the discretization error. The theoretical results are illustrated by numerical computations.

**Key words.** obstacle problems, adaptive finite element methods, multilevel preconditioning, a posteriori error estimates

**AMS subject classifications.** 65N30, 65N50, 65N55, 35J85, 49J40

**1. Introduction.** Given a closed subspace  $V \subset H^1(\Omega)$ ,  $\Omega$  being a bounded polygonal domain in the Euclidean space  $\mathbb{R}^2$ , we consider obstacle problems of the form

$$(1) \quad \text{find } u \in K \text{ such that } \mathcal{J}(u) \leq \mathcal{J}(v), \quad v \in K,$$

for the energy functional  $\mathcal{J}$ ,

$$\mathcal{J}(v) = \frac{1}{2}a(v, v) - \ell(v), \quad v \in V,$$

and a closed, convex set  $K \subset V$ ,

$$K = \{v \in V \mid v(x) \leq \varphi(x) \text{ a.e. in } \Omega\}.$$

Assuming that  $\mathcal{J}$  is induced by a symmetric  $V$ -elliptic bilinear form  $a(\cdot, \cdot)$ ,

$$a(v, w) = \int_{\Omega} \sum_{i,j=1}^2 a_{ij} \partial_i v \partial_j w \, dx,$$

and some functional  $\ell \in V'$ , it is well-known that (1) is equivalent to the variational inequality

$$(2) \quad \text{find } u \in K \text{ such that } a(u, u - v) \leq \ell(u - v), \quad v \in K.$$

For the sake of simplicity we restrict our considerations to the case  $V = H_0^1(\Omega)$ . To ensure existence and uniqueness of the solution  $u$  of (1) and (2), respectively, we assume  $\varphi \in H^1(\Omega)$ ,  $\varphi \geq 0$  almost everywhere on  $\Gamma = \partial\Omega$ , and  $a_{ij} \in L^\infty(\Omega)$  satisfying

$$(a) \quad a_{ij}(x) = a_{ji}(x), \quad 1 \leq i, j \leq 2,$$

$$(3) \quad (b) \quad \alpha_0 |\xi|^2 \leq \sum_{i,j=1}^2 a_{ij}(x) \xi_i \xi_j \leq \alpha_1 |\xi|^2, \quad \xi \in \mathbb{R}^2, \quad 0 < \alpha_0 \leq \alpha_1$$

\* Received by the editors July 9, 1992; accepted for publication (in revised form) February 16, 1993.

† Mathematisches Institut der Technischen Universität München Arcisstrasse 21, D-80333 München, Germany (rohopp@mathematik.tu-muenchen.de).

‡ Konrad-Zuse-Zentrum für Informationstechnik Berlin Heilbronner Strasse 10, D-10711 Berlin, Germany (kornhuber@sc.ZIB-Berlin.DE).

for almost all  $x \in \Omega$ .

Discretizing (2) in space by continuous, piecewise linear finite elements with respect to a triangulation of  $\Omega$ , standard numerical schemes for the solution of the resulting finite-dimensional variational inequality are projected relaxation methods (e.g., [17]). These iterative methods typically suffer from rapidly deteriorating convergence rates when proceeding to more and more refined triangulations, which renders them inefficient from a numerical point of view. However, this drawback can be overcome by using multilevel techniques with respect to a hierarchy of triangulations. Multigrid approaches to obstacle problems have been developed by various authors ([10], [18]–[21], [30], [31]). For obstacle-type problems an alternative to projected relaxation is to use some sort of linearization techniques based on active set strategies (e.g., [18]–[20]). This is an iterative scheme in which at each iteration step a set of active constraints is prespecified and then a linear subproblem has to be solved for the computation of the new iterate. Note that the multigrid techniques used in [18]–[20] consist of outer and inner iterations where the outer iteration is an active set strategy and the inner iterations are multigrid iterations for the approximate solution of the auxiliary problems.

Because the coefficient matrices of the auxiliary systems are symmetric positive definite for the obstacle problems under consideration, an alternative choice for the inner iterations are preconditioned conjugate gradient (pcg) methods, especially those based on multilevel preconditioners such as Yserentant's hierarchical basis preconditioner [38] or the BPX-preconditioner [9]. A related approach has been proposed by Schwenkert [36] in which relaxation methods have been applied with respect to hierarchical bases.

For the adaptive construction of a suitable hierarchy of triangulations efficient and reliable a posteriori error estimates are required. While a variety of well-established results are available in the case of linear elliptic problems (see [3], [14], [23], [37] for further references) the situation is less clear in the case of obstacle problems. Recently, the concepts introduced in [14] were extended and applied successfully to a special obstacle problem arising in semiconductor device simulation [26]. A more detailed investigation of this approach will be one of the subjects of this paper. A posteriori error estimates for the penalty method together with strategies for the adaptive choice of a space-dependent penalty parameter and the meshsize were given in [24].

The paper is organized as follows. After a brief discussion of the active-set strategy proposed in [19], we will focus on the construction and analysis of multilevel preconditioners providing the efficient solution of the arising linear subproblems. In particular, we will derive two variants of hierarchical basis type by suitable modifications of the standard hierarchical basis preconditioner. It will turn out that both variants are performing asymptotically as in the unconstrained case, but that only one of them is robust with respect to the regularity of the free boundary. Inspired by a paper of Dryja and Widlund [15], the preconditioners will be regarded as multilevel additive Schwarz (MAS) methods. This abstract framework allows for obvious extensions to other variants of the MAS method, in particular to the BPX-preconditioner. By comparing the actual approximation with another approximation of higher accuracy, we will derive semilocal and local a posteriori error estimates, followed by a detailed analysis of their efficiency and reliability. The final chapter is devoted to some numerical experiments supporting the theoretical findings.

**2. Outer-inner iterations.** Let  $\mathcal{T}$  denote a triangulation of the computational domain  $\Omega \subset \mathbb{R}^2$ . We assume that  $\mathcal{T}$  is regular in the sense that the intersection of

two triangles  $t, t' \in \mathcal{T}$  contains a common edge, a common vertex, or is empty. The sets of vertices  $p$  and edges  $e$  that are not part of the boundary  $\partial\Omega$  are called  $\mathcal{N}$  and  $\mathcal{E}$ , respectively. We approximate  $V$  by the subspace  $\mathcal{S}$  of continuous, piecewise linear finite elements vanishing on the boundary  $\partial\Omega$  with the associated nodal basis  $\lambda_p, p \in \mathcal{N}$ , of  $\mathcal{S}$  defined by  $\lambda_p(q) = \delta_{pq}, p, q \in \mathcal{N}$  (Kronecker delta).

Furthermore, let  $\varphi_{\mathcal{T}} \in \mathcal{S}$  be a discrete obstacle approximating the given obstacle  $\varphi$  in an appropriate sense. For example,  $\varphi_{\mathcal{T}}$  may be chosen as the  $L^2$ -projection of  $\varphi$  onto  $\mathcal{S}$  or, if  $\varphi \in C(\bar{\Omega})$ , as the  $\mathcal{S}$ -interpolate. Correspondingly, we denote by  $K_{\mathcal{T}} = \{v \in \mathcal{S} | v \leq \varphi_{\mathcal{T}}\}$  the sets of discrete constraints. Then the finite element approximation of (1) amounts to the computation of an element  $u_{\mathcal{T}} \in K_{\mathcal{T}}$  satisfying

$$(4) \quad a(u_{\mathcal{T}}, u_{\mathcal{T}} - v) \leq \ell(u_{\mathcal{T}} - v), \quad v \in K_{\mathcal{T}}.$$

It is easy to see that the finite-dimensional variational inequality (4) is equivalent to a linear complementarity problem [13].

LEMMA 2.1. *An element  $u_{\mathcal{T}} \in K_{\mathcal{T}}$  is a solution to (4) if and only if the vector  $\underline{u} \in \mathbb{R}^N, N := |\mathcal{N}|$  with components  $u_p = u_{\mathcal{T}}(p), p \in \mathcal{N}$ , satisfies*

$$(5) \quad \max(A\underline{u} - \underline{b}, \underline{u} - \underline{\varphi}) = 0,$$

where  $A$  is the  $N \times N$  stiffness matrix with entries  $a_{pq} = a(\lambda_q, \lambda_p), p, q \in \mathcal{N}$ , and  $\underline{b} \in \mathbb{R}^N$  and  $\underline{\varphi} \in \mathbb{R}^N$  are the vectors with components  $b_p = \ell(\lambda_p)$  and  $\varphi_p = \varphi_{\mathcal{T}}(p), p \in \mathcal{N}$ . Note that (5) has to be understood in terms of its components.

*Proof.* Let  $u_{\mathcal{T}} \in K_{\mathcal{T}}$  be the solution of (4). Then  $A\underline{u} \leq \underline{b}$ , which can be deduced by choosing  $v = u_{\mathcal{T}} - z$  in (4) with arbitrarily given  $z \in \mathcal{S}, z \geq 0$ . Since  $\underline{u} \leq \underline{\varphi}$ , we thus have  $(\underline{u} - \underline{\varphi})^T (A\underline{u} - \underline{b}) \geq 0$ . But  $v = \varphi_{\mathcal{T}}$  in (4) gives  $(\underline{u} - \underline{\varphi})^T (A\underline{u} - \underline{b}) \leq 0$ , hence  $(\underline{u} - \underline{\varphi})^T (A\underline{u} - \underline{b}) = 0$ , proving (5). The converse statement is obvious.  $\square$

In the following algorithm we will consider an outer-inner iteration technique for the numerical solution of the complementarity problem (5). The outer iterations are governed by an active-set strategy as presented in [19], [20]:

*Outer iteration.* (active-set strategy):

Step 1. Choose a startvector  $\underline{u}^{(0)} \in \mathbb{R}^N$ .

Step 2. Given  $\underline{u}^{(\nu)} \in \mathbb{R}^N, \nu \geq 0$ , determine  $\mathcal{N}^\bullet \subset \mathcal{N}$  as the set of points  $p \in \mathcal{N}$  such that  $(A\underline{u}^{(\nu)} - \underline{b})_p \leq (\underline{u}^{(\nu)} - \underline{\varphi})_p$  and set  $\mathcal{N}^\circ := \mathcal{N} \setminus \mathcal{N}^\bullet$ . Then compute  $\underline{u}^{(\nu+1)} \in \mathbb{R}^N$  from the splitting

$$(6) \quad \underline{u}^{(\nu+1)} = \underline{u}^\bullet + \underline{u}^\circ,$$

where

$$(7) \quad u_p^\bullet = \varphi_p, p \in \mathcal{N}^\bullet, \quad u_p^\bullet = 0, p \in \mathcal{N}^\circ$$

and  $\underline{u}^\circ$  satisfies

$$(8) \quad u_p^\circ = 0, \quad p \in \mathcal{N}^\bullet$$

and

$$(9) \quad A\underline{u}^\circ = \underline{b} - A\underline{u}^\bullet.$$

It is obvious that the computation of the iterate  $\underline{u}^{(\nu+1)}$  according to (9) actually requires the solution of a ‘‘reduced,’’ i.e., lower-dimensional linear system.

The set  $\mathcal{N}^\bullet$  is called active, since in view of  $u_p^{(\nu+1)} = \varphi_p$ ,  $p \in \mathcal{N}^\bullet$ , it contains the nodal points where the obstacle is active. Correspondingly,  $\mathcal{N}^\circ$  is said to be the inactive set. Introducing a corresponding splitting of the finite element space  $\mathcal{S} = \mathcal{S}^\circ \oplus \mathcal{S}^\bullet$  in linear subspaces  $\mathcal{S}^\circ$ ,  $\mathcal{S}^\bullet \subset \mathcal{S}$  defined by

$$(10) \quad \mathcal{S}^\circ = \{v \in \mathcal{S} \mid v(p) = 0, p \in \mathcal{N}^\bullet\}, \quad \mathcal{S}^\bullet = \{v \in \mathcal{S} \mid v(p) = 0, p \in \mathcal{N}^\circ\},$$

the reduced system (9) can be rewritten as the variational equality

$$(11) \quad \text{find } u^\circ \in \mathcal{S}^\circ \text{ such that } a(u^\circ, v) = \ell(v) - a(u^\bullet, v), \quad v \in \mathcal{S}^\circ,$$

with solution  $u^\circ \in \mathcal{S}^\circ$ , where  $u^\bullet \in \mathcal{S}^\bullet$  is defined by  $u^\bullet(p) = u_p^\bullet$ .

*Remark 2.1.* If (9) (respectively, (11)) is solved exactly and  $A$  is an M-matrix, it can be shown that for an arbitrarily given initial iterate  $\underline{u}^{(0)}$  the sequence  $\underline{u}^{(\nu)}$ ,  $\nu \geq 1$ , of iterates is monotonically decreasing and converging to the unique solution  $\underline{u}$  of (5). Moreover, there is numerical evidence that the approximate solution of the linear subproblems up to some accuracy  $\kappa_0$  provides satisfying results as soon as  $\kappa_0$  is chosen small enough. See [19] and [20] for details. In the inexact case, the convergence of a related most-constrained strategy has been proved in [18], providing a stopping criterion for the inner iteration. However, this strategy turns out to be much too pessimistic in actual computation, leading to a prohibitively large number of outer iteration steps.

In contrast to [19] and [20] (where multigrid techniques were used), this paper focuses on multilevel preconditioned conjugate gradient (cg) iterations that for well-known reasons are suited to be used within an adaptive finite-element code (FEM). For an introduction to the preconditioned cg method we refer to [1], while the construction of appropriate multilevel preconditioners will be subject of the next chapter.

**3. Additive Schwarz methods and hierarchical bases.** Let  $\mathcal{T}_0$  be an intentionally coarse regular triangulation of  $\Omega$ .

The triangulation  $\mathcal{T}_0$  is refined several times, providing a sequence of triangulations  $\mathcal{T}_0, \mathcal{T}_1, \dots, \mathcal{T}_j$  and a corresponding sequence of nested finite element spaces  $\mathcal{S}_0 \subset \mathcal{S}_1 \subset \dots \subset \mathcal{S}_j$ . The underlying refinement process described in the sequel is standard in the literature on multilevel preconditioning [3]–[6], [8], [14], [39]. Note that this refinement in general does not coincide with the actual refinement process performed by some finite element code. Nevertheless, the triangulations  $\mathcal{T}_0, \mathcal{T}_1, \dots, \mathcal{T}_j$  are available without any computational effort, if the underlying data structures are chosen properly [3], [28], [33], [34].

A triangle  $t \in \mathcal{T}_k$  is refined either by subdividing it into four congruent subtriangles or by connecting one of its vertices with the midpoint of the opposite side. The first case is called regular (red) refinement and the resulting triangles are regular, as are the triangles of the initial triangulation  $\mathcal{T}_0$ . The second case is called irregular (green) refinement and results in two irregular triangles. Because new points should be generated only by regular refinement, we introduce the following rule:

(T1) Each vertex of  $\mathcal{T}_{k+1}$  that does not belong to  $\mathcal{T}_k$  is a vertex of a regular triangle.

Note that irregular refinement is potentially dangerous because the interior angles are reduced. Hence we add the following rule:

(T2) Irregular triangles must not be further refined.

We say that a refined triangle is the father of the resulting triangles, which in turn are called sons. We define the depth of a given triangle  $t \in \bigcup_{k=0}^j \mathcal{T}_k$  as the number of

ancestors of  $t$ . Of course, the depth of all triangles  $t \in \mathcal{T}_k$  is bounded by  $k$ . We have the final rule:

(T3) Only triangles  $t \in \mathcal{T}_k$  of depth  $k$  may be refined for the construction of  $\mathcal{T}_{k+1}$ ,  $0 \leq k \leq j$ .

As a consequence of (T3), the whole sequence  $\mathcal{T}_0, \mathcal{T}_1, \dots, \mathcal{T}_j$  can be uniquely reconstructed from the initial triangulation  $\mathcal{T}_0$  and the final triangulation  $\mathcal{T}_j$  alone, neglecting the preceding dynamic refinement process. Recall that in actual computations we may choose the data structures representing the triangulations cleverly so that the sequence  $\mathcal{T}_0, \mathcal{T}_1, \dots, \mathcal{T}_j$  is explicitly given. Note that the subscript  $j$  in general does *not* coincide with the number of refinement steps,  $l \geq j$ , which were necessary to create  $\mathcal{T}_j$  from  $\mathcal{T}_0$  by the actual finite element code. In practical calculations the difference  $l - j$  of the refinement level  $l$  and the maximal depth  $j$  can be used to judge the quality of the implemented refinement strategy.

Of course, adaptive refinement should be based on reliable a posteriori error estimates, which will be considered in the following chapter. For the moment let us assume that a hierarchy  $\mathcal{T}_0, \mathcal{T}_1, \dots, \mathcal{T}_j$  with the property (T1–T3) is available. We further assume that we have a disjoint splitting  $\mathcal{N}_j = \mathcal{N}_j^\bullet \cup \mathcal{N}_j^\circ$ , which may result from an active set strategy applied to (4) with respect to the triangulation  $\mathcal{T} = \mathcal{T}_j$ . Recall that this splitting is supposed to change in each outer iteration step. In the sequel we will deal with the construction of two multilevel preconditioners of hierarchical basis type to provide an efficient iterative solution of the corresponding reduced system:

$$(12) \quad \text{find } u_j^\circ \in \mathcal{S}_j^\circ \text{ such that } a(u_j^\circ, v) = \ell(v) - a(u_j^\bullet, v), \quad v \in \mathcal{S}_j^\circ.$$

For this purpose we provide a decomposition  $\mathcal{N}_k = \mathcal{N}_k^\bullet \cup \mathcal{N}_k^\circ$  of the sets  $\mathcal{N}_k$  of the nodal points on the lower levels  $0 \leq k < j$  by means of the definition

$$(13) \quad \mathcal{N}_k^\bullet = \mathcal{N}_k \cap \mathcal{N}_j^\bullet, \quad \mathcal{N}_k^\circ = \mathcal{N}_k \setminus \mathcal{N}_k^\bullet, \quad 0 \leq k \leq j - 1.$$

For  $0 \leq k \leq j$  and  $p \in \mathcal{N}_k$  we refer to  $\lambda_p^{(k)} \in \mathcal{S}_k$  as the level  $k$  nodal basis function having  $p$  as its supporting point, i.e.,  $\lambda_p^{(k)}(p) = 1$ . The well-known hierarchical basis of the whole space  $\mathcal{S}_j$  (cf. Yserentant [38]) is given by

$$\Lambda_0 = \{\lambda_p^{(0)} \mid p \in \mathcal{N}_0\}, \quad \Lambda_k = \{\lambda_p^{(k)} \mid p \in \mathcal{N}_k \setminus \mathcal{N}_{k-1}\}, \quad 1 \leq k \leq j,$$

denoting  $\Lambda = \bigcup_{k=1}^j \Lambda_k$ . According to (10) the splitting (13) induces the subspaces  $\mathcal{S}_k^\circ = \text{span}\{\lambda_p^{(k)} \mid p \in \mathcal{N}_k^\circ\} \subset \mathcal{S}_k$ ,  $0 \leq k \leq j$ . Collecting the hierarchical basis functions with inactive supporting points according to

$$(14) \quad \hat{\Lambda}_0 := \{\lambda_p^{(0)} \mid p \in \mathcal{N}_0^\circ\}, \quad \hat{\Lambda}_k := \{\lambda_p^{(k)} \mid p \in \mathcal{N}_k^\circ \setminus \mathcal{N}_{k-1}^\circ\}, \quad 1 \leq k \leq j,$$

we denote  $\hat{\Lambda}_H = \bigcup_{k=1}^j \hat{\Lambda}_k$ . However, the hierarchical decomposition of functions  $v \in \mathcal{S}_j^\circ$  cannot be given in the standard way, since the subsets  $\hat{\Lambda}_0$  and  $\hat{\Lambda}_k$  of  $\mathcal{S}_j$  in general are not contained in  $\mathcal{S}_j^\circ$ . This is due to the fact that functions  $v \in \mathcal{S}_{k-1}^\circ$ ,  $1 \leq k \leq j$ , in general do not vanish in active nodal points  $p \in \mathcal{N}_k^\bullet \setminus \mathcal{N}_{k-1}^\bullet$  appearing on the subsequent level  $k$ . We will modify such functions by means of suitable truncation operators  $T_k : \mathcal{S}_j \rightarrow \mathcal{S}_k^\circ$ ,  $0 \leq k \leq j$ , defined by

$$(15) \quad T_k v = \sum_{p \in \mathcal{N}_k^\circ} v(p) \lambda_p^{(k)}.$$

Note that  $T_k v = v$ ,  $v \in \mathcal{S}_k^\circ$ . Now a feasible multilevel splitting of  $\mathcal{S}_j^\circ$  is defined by successive truncation of the standard hierarchical basis elements

$$(16) \quad \Lambda_k^{(1)} = T_{j,k} \hat{\Lambda}_k, \quad T_{j,k} = T_j \dots T_k, \quad 0 \leq k \leq j.$$

We will consider a second multilevel splitting, which is based on a more restrictive choice of coarse grid functions. For this reason we define

$$(17) \quad \mathcal{N}_k^{\circ, \text{reg}} = \{p \in \mathcal{N}_k^\circ \mid T_{j,k} \lambda_p^{(k)} = \lambda_p^{(k)}\}, \quad 0 \leq k \leq j$$

and  $\mathcal{N}_k^{\bullet, \text{reg}} = \mathcal{N}_k \setminus \mathcal{N}_k^{\circ, \text{reg}}$ . Obviously, we have  $\mathcal{N}_j^{\circ, \text{reg}} = \mathcal{N}_j^\circ$ . It is easily seen by induction that for  $0 < k \leq j$  the set  $\mathcal{N}_{k-1}^{\circ, \text{reg}}$  consists of all  $p \in \mathcal{N}_{k-1} \cap \mathcal{N}_k^{\circ, \text{reg}}$  whose  $k$ -neighbors  $q \in \mathcal{N}_k \setminus \mathcal{N}_{k-1}$  are also contained in  $\mathcal{N}_k^{\circ, \text{reg}}$ . As usual,  $p, q \in \mathcal{N}_k$  are called  $k$ -neighbors if there is an edge  $e = (p, q) \in \mathcal{E}_k$ . Now the standard hierarchical splitting with respect to  $\mathcal{N}_k^{\circ, \text{reg}}$ ,  $0 \leq k \leq j$ , is given by

$$(18) \quad \Lambda_0^{(2)} = \{\lambda_p^{(0)} \mid p \in \mathcal{N}_0^{\circ, \text{reg}}\}, \quad \Lambda_k^{(2)} = \{\lambda_p^{(k)} \mid p \in \mathcal{N}_k^{\circ, \text{reg}} \setminus \mathcal{N}_{k-1}^{\circ, \text{reg}}\}, \quad 1 \leq k \leq j.$$

Note that a restriction of the active set, which is similar to (17), was used in [19]. In the context of hierarchical bases, (17) was proposed by Yserentant [40].

*Remark 3.1.* The difference between  $\Lambda_k^{(1)}$  and  $\Lambda_k^{(2)}$  is illustrated in Fig. 1, where for ease of exposition we have considered the one-dimensional case.

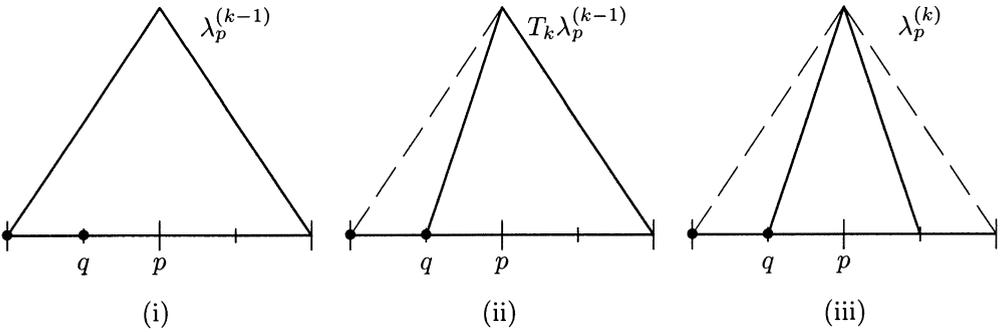


FIG. 1.

In particular, Fig. 1 (i) represents a level  $k - 1$  basis function  $\lambda_p^{(k-1)}$ , with supporting point  $p \in \mathcal{N}_{k-1}^\circ \setminus \mathcal{N}_{k-1}^{\circ, \text{reg}}$  having a level  $k$  active neighbor  $q \in \mathcal{N}_k^\bullet$  on the left. Figures 1 (ii) and (iii) display the basis functions  $T_k \lambda_p^{(k-1)}$  and  $\lambda_p^{(k)}$  selected in (16) and (18), respectively. Note that  $T_k \lambda_p^{(k-1)}$  generally results in a “nonsymmetric” truncation, while the choice of the higher level basis function  $\lambda_p^{(k)}$  may be regarded as a “symmetric” cut.

As proposed in [15], the hierarchical basis preconditioners obtained from (16) and (18) will be treated in the framework of additive Schwarz methods. For recent results on the BPX preconditioner as an additive Schwarz method, we refer the reader to Bornemann [6] and Zhang [41]. Because the following definitions and assertions do not differ for  $\Lambda_k^{(\mu)}$ ,  $\mu = 1, 2$ , the index  $\mu$  is skipped for notational convenience.

Let  $V_0^{(\mu)} = \text{span } \Lambda_0^{(\mu)}$  and  $V_\lambda^{(\mu)} = \text{span}\{\lambda\}$ ,  $\lambda \in \Lambda_H^{(\mu)} = \bigcup_{k=1}^j \Lambda_k^{(\mu)}$  for  $\mu = 1, 2$ . The direct subspace decomposition

$$(19) \quad \mathcal{S}_j^\circ = V_0 \oplus \bigoplus_{\lambda \in \Lambda_H} V_\lambda$$

of  $\mathcal{S}_j^\circ$  gives rise to an additive Schwarz method. This in turn provides the following reformulation of the original problem (12):

$$Pu_j^\circ = \ell',$$

where

$$P = P_0 + \sum_{\lambda \in \Lambda_H} P_\lambda$$

is the sum of the Ritz projections  $P_0 : \mathcal{S}_j^\circ \rightarrow V_0$ ,  $P_\lambda : \mathcal{S}_j^\circ \rightarrow V_\lambda$ ,  $\lambda \in \Lambda_H$ , defined by

$$a(P_\nu w, v) = a(w, v), \quad v \in V_\nu, \quad \nu = 0, \lambda,$$

for each  $w \in \mathcal{S}_j^\circ$ , and  $\ell' \in (\mathcal{S}_j^\circ)'$  is chosen appropriately. See, for example, [16] for details. Denoting by  $(\cdot, \cdot)$  the standard  $L^2$  inner product, we introduce the  $L^2$  projections  $Q_0 : \mathcal{S}_j^\circ \rightarrow V_0$ ,  $Q_\lambda : \mathcal{S}_j^\circ \rightarrow V_\lambda$  and the representation operators  $A_0 : V_0 \rightarrow V_0$ ,  $A_\lambda : V_\lambda \rightarrow V_\lambda$ ,  $\lambda \in \Lambda_H$  defined by

$$(Q_\nu w, v) = (w, v), \quad v \in V_\nu,$$

for each  $w \in \mathcal{S}_j^\circ$  and

$$(A_\nu w, v) = a(w, v), \quad v \in V_\nu,$$

for each  $w \in V_\nu$ ,  $\nu = 0, \lambda$ . Since  $A_\nu P_\nu = Q_\nu A$ ,  $\nu = 0, \lambda$ , the operator  $P$  may be rewritten as

$$P = H_j A_j,$$

where  $H_j$  stands for the preconditioner

$$H_j = A_0^{-1} Q_0 + \sum_{\lambda \in \Lambda_H} A_\lambda^{-1} Q_\lambda,$$

and  $A_j$  is the representation operator of  $a(\cdot, \cdot)$  on  $\mathcal{S}_j^\circ \times \mathcal{S}_j^\circ$ . Evaluation of  $A_\lambda^{-1} Q_\lambda$  leads to

$$(20) \quad H_j^{(\mu)} = (A_0^{(\mu)})^{-1} Q_0^{(\mu)} + \sum_{\lambda \in \Lambda_H^{(\mu)}} \frac{(\cdot, \lambda)}{a(\lambda, \lambda)} \lambda, \quad \mu = 1, 2.$$

In light of Remark 3.1, we will refer to  $H_j^{(1)}$  and its variants as the “nonsymmetric” preconditioners and to  $H_j^{(2)}$  as the “symmetric” preconditioner, respectively.

Let us briefly discuss some modifications of the preconditioners  $H_j^{(\mu)}$ ,  $\mu = 1, 2$ . The evaluation of  $(A_0^{(\mu)})^{-1} Q_0^{(\mu)}$  requires the solution of a linear system for the stiffness matrix given by  $a(\cdot, \cdot)$  restricted to  $V_0^{(\mu)} \times V_0^{(\mu)}$ ,  $\mu = 1, 2$ . Due to the definition (16), the entries of  $A_0^{(1)}$  and  $a(\lambda, \lambda)$ ,  $\lambda \in \Lambda_H^{(1)}$  may change with each step of the outer iteration. To avoid the corresponding evaluations of the quadratic form  $a(\cdot, \cdot)$  the preconditioner  $H_j^{(1)}$  may be replaced by

$$(21) \quad \hat{H}_j^{(1)} = T_{j,0} \hat{A}_0^{-1} \hat{Q}_0 + \sum_{k=1}^j \sum_{\lambda \in \hat{\Lambda}_k} \frac{(\cdot, T_{j,k} \lambda)}{a(\lambda, \lambda)} T_{j,k} \lambda,$$

where  $\hat{A}_0$  is the representation of  $a(\cdot, \cdot)$  restricted to  $\mathcal{S}_0^\circ \times \mathcal{S}_0^\circ$  and  $\hat{Q}_0$  denotes the  $L^2$  projection to  $\mathcal{S}_0^\circ$ , respectively. Note that a related modification of  $H_j^{(2)}$  is not necessary, as only the selection and not the shape of the involved hierarchical basis functions is depending on the actual active set  $\mathcal{N}_j^\bullet$ . Still, the linear system on the coarsest level is supposed to change with each outer iteration step, each time causing a Cholesky decomposition of the new coefficient matrix. To reduce the computational effort, we may replace the matrix by its diagonal or even by the identity matrix (see [38] for a further discussion). In the case of rapidly varying coefficients, frequently occurring in practical problems, the jumps should be incorporated in the preconditioners. We refer to Yserentant [39] for details.

Note that existing implementations of the standard hierarchical basis preconditioner are easily changed to (21) by simply neglecting the contributions from active points [22]. For a similar application of truncated hierarchical basis functions to obstacle problems, we refer to [36].

The final part of this section will provide condition number estimates for both the nonsymmetric and the symmetric cases. The subsequent analysis will be guided by the following lemma on abstract additive Schwarz methods.

LEMMA 3.1. (i) Assume that for all  $v \in \mathcal{S}_j^\circ$  there is a splitting  $v = v_0 + \sum_{\lambda \in \Lambda_H} v_\lambda$  such that

$$(22) \quad c \left\{ a(v_0, v_0) + \sum_{\lambda \in \Lambda_H} a(v_\lambda, v_\lambda) \right\} \leq a(v, v)$$

holds for some fixed positive constant  $c$ . Then we have the estimate

$$ca(v, v) \leq a(Pv, v), \quad v \in \mathcal{S}_j^\circ.$$

(ii) Assume that for all splittings  $v = v_0 + \sum_{\lambda \in \Lambda_H} v_\lambda$  of  $v \in \mathcal{S}_j^\circ$  the estimate

$$(23) \quad a(v, v) \leq C \left\{ a(v_0, v_0) + \sum_{\lambda \in \Lambda_H} a(v_\lambda, v_\lambda) \right\}$$

holds for some fixed positive constant  $C$ . Then we have the estimate

$$a(Pv, v) \leq Ca(v, v), \quad v \in \mathcal{S}_j^\circ.$$

*Proof.* The assertion (i) is the well-known lemma of Lions [29]. To prove the second assertion, we apply (23) to the splitting  $Pv = P_0v + \sum_{\lambda \in \Lambda_H} P_\lambda v$  for some fixed  $v \in \mathcal{S}_j^\circ$  to obtain

$$a(Pv, Pv) \leq C \left\{ a(P_0v, P_0v) + \sum_{\lambda \in \Lambda_H} a(P_\lambda v, P_\lambda v) \right\} = Ca(Pv, v),$$

which completes the proof.  $\square$

*Remark 3.2.* The assumptions (22) and (23) can be regarded as an asymptotic orthogonality of the subspaces  $V_0, V_\lambda, \lambda \in \Lambda_H$ . Note that (23) is frequently established by strengthened Cauchy–Schwarz inequalities measuring the angles between  $V_0, V_\lambda, \lambda \in \Lambda_H$  with respect to  $a(\cdot, \cdot)$  or any other symmetric bilinear form that generates a uniformly equivalent norm on  $\mathcal{S}_j$ . We will use this approach later on.

In addition to the usual (semi-) norms  $\|\cdot\|_0$  and  $|\cdot|_1$  of  $L^2(\Omega)$  and  $H^1(\Omega)$ , we will make use of the semi-inner product

$$(v, w)_{1, \Omega_0} = \sum_{i=1}^2 \int_{\Omega_0} \partial_i v \partial_i w \, dx, \quad v, w \in H^1(\Omega_0)$$

for measurable  $\Omega_0 \subset \Omega$  with the induced seminorm  $|v|_{1, \Omega_0} = (v, v)_{1, \Omega_0}^{1/2}$ . We continue by introducing the interpolation operators  $I_k : \mathcal{S}_j \rightarrow \mathcal{S}_k$  by

$$I_k v = \sum_{p \in \mathcal{N}_k} v(p) \lambda_p^{(k)}, \quad 0 \leq k \leq j.$$

Finally, constants depending only on the ellipticity (3) and the shape regularity of  $\mathcal{T}_0$  will be denoted by  $c$  or  $C$ . Other parameters will be indicated explicitly.

We take up the analysis of the preconditioners with the following technical lemma.

**LEMMA 3.2.** *For some fixed  $k$ ,  $0 < k \leq j$ , let  $\mathcal{S}_k^* \subset \mathcal{S}_k$  denote the subspace of all  $v_k \in \mathcal{S}_k$  vanishing in  $\mathcal{N}_k^* \subset \mathcal{N}_k$ . Assume that each point  $p \in \mathcal{N}_{k-1}$  is either contained in  $\mathcal{N}_k^*$  or has at least one  $k$ -neighbor  $q \in \mathcal{N}_k \setminus \mathcal{N}_{k-1}$  that is contained in  $\mathcal{N}_k^*$ . Then we have the estimate*

$$\sum_{p \in \mathcal{N}_k} |v_k(p) \lambda_p^{(k)}|_1^2 \leq C |v_k|_1^2, \quad v_k \in \mathcal{S}_k^*.$$

*Proof.* Let  $p \in \mathcal{N}_{k-1}$ . Then from the assumptions on  $\mathcal{N}_{k-1}$  and  $\mathcal{N}_k^*$ , the seminorm  $|\cdot|_1$  is a norm on the restriction of  $v_k \in \mathcal{S}_k^*$  to  $\text{supp } \lambda_p^{(k-1)}$ . As norms on finite-dimensional spaces are equivalent, we immediately have

$$\sum_{q \in \mathcal{N}_k \cap \text{supp } \lambda_p^{(k-1)}} |v_k(q) \lambda_q^{(k)}|_1^2 \leq c_{k,p} |v_k|_{1, \text{supp } \lambda_p^{(k-1)}}^2, \quad v_k \in \mathcal{S}_k^*,$$

and as a consequence of the uniform shape regularity of  $\mathcal{T}_k$  we obtain  $c_{k,p} \leq c$  uniformly in  $k$  and  $p$ . Summing up over all  $p \in \mathcal{N}_{k-1}$  gives the assertion.  $\square$

Crucial to the analysis of the nonsymmetric preconditioners  $H_j^{(1)}$  and  $\hat{H}_j^{(1)}$  is the following assumption on the splitting  $\mathcal{N}_j = \mathcal{N}_j^\bullet \cup \mathcal{N}_j^\circ$ :

(R) There is a nonnegative constant  $k_0$  independent of  $k$  such that

$$(24) \quad T_{j,k} \lambda = T_{k+k_0,k} \lambda, \quad \lambda \in \hat{\Lambda}_k, \quad j \geq k + k_0.$$

*Remark 3.3.* The condition (R) states that subsequent truncation of level  $k$  functions uniformly becomes stationary after  $k_0$  steps. From heuristic arguments we can expect that for each level  $k$  we can find a number  $k_0 = k_0(k)$  satisfying (24) if the free boundary is a lower-dimensional manifold that is properly approximated by the finite element discretization and the underlying active set strategy. It is additionally required by the condition (R) that these numbers be uniformly bounded. The condition (R) will typically be applied as in the proof of the following lemma.

**LEMMA 3.3.** *Assume that (R) is satisfied. Then there exist constants  $c(k_0)$ ,  $C(k_0)$  such that*

$$(25) \quad c(k_0) a(v, v) \leq a(T_{j,k} v, T_{j,k} v) \leq C(k_0) a(v, v)$$

holds for all  $v \in \text{span } \hat{\Lambda}_k$ ,  $k = 1, \dots, j$ .

*Proof.* Without loss of generality, we assume  $j \geq k + k_0$ . Let  $v \in \text{span } \hat{\Lambda}_k$ . Then the condition (R) provides

$$(26) \quad T_{j,k}v = T_{k+k_0,k}v \in \mathcal{S}_{k+k_0}.$$

Let  $t$  be a triangle of  $\mathcal{T}_{k-1}$ . As  $v$  vanishes in the vertices of  $t$  and the space of functions of  $\mathcal{S}_{k+k_0}$  restricted to  $t$  is finite dimensional, there are positive constants  $c(k_0)$  and  $C(k_0)$  depending on  $k_0$  and a lower bound for the interior angles with

$$(27) \quad c(k_0) |v|_{1,t} \leq |T_{j,k}v|_{1,t} \leq C(k_0) |v|_{1,t}, \quad v \in \text{span } \hat{\Lambda}_k.$$

Summing up over  $t \in \mathcal{T}_{k-1}$ , the assertion follows from the ellipticity of  $a(\cdot, \cdot)$ .  $\square$

We are ready to establish lower and upper bounds for the nonsymmetric preconditioners  $H_j = H_j^{(1)}, \hat{H}_j^{(1)}$ .

**THEOREM 3.4.** *Assume that the regularity condition (R) holds. Then there exist constants  $K_0, K_1$  depending only on  $\alpha_0, \alpha_1$  in (3), the shape regularity of  $\mathcal{T}_0$  and the constant  $k_0$  in (R) such that the estimate*

$$K_0(j+1)^{-2}a(v, v) \leq a(H_j A_j v, v) \leq K_1 a(v, v), \quad H_j = H_j^{(1)}, \hat{H}_j^{(1)}$$

holds for all  $v \in \mathcal{S}_j^\circ$ .

*Proof.* Let us first consider the case  $H_j = H_j^{(1)}$ . To verify the assumption of Lemma 3.1(i), we consider the splitting

$$(28) \quad v = \tilde{v}_0 + \sum_{\tilde{\lambda} \in \Lambda_H^{(1)}} v_{\tilde{\lambda}}, \quad \tilde{v}_0 \in V_0^{(1)}, \quad v_{\tilde{\lambda}} = v(\tilde{\lambda})\tilde{\lambda}, \quad \tilde{\lambda} \in \Lambda_H^{(1)}$$

of some fixed  $v \in \mathcal{S}_j^\circ$ . As (19) provides a direct splitting of  $\mathcal{S}_j^\circ$ , this representation is unique. However, there is another decomposition of  $v$  with respect to the standard hierarchical basis  $\Lambda$  of  $\mathcal{S}_j$ ,

$$v = v_0 + \sum_{\lambda \in \Lambda} v_\lambda, \quad v_0 \in \mathcal{S}_0, \quad v_\lambda = v(\lambda)\lambda, \quad \lambda \in \Lambda.$$

Observe that we have

$$\tilde{v}_0 = T_{j,0}v_0, \quad v_{\tilde{\lambda}} = T_{j,k}v_\lambda, \quad \lambda \in \hat{\Lambda}_k.$$

Together with Lemma 3.3 and the equivalence of norms on the coarse space  $V_0^{(1)}$ , this gives

$$(29) \quad a(\tilde{v}_0, \tilde{v}_0) + \sum_{\tilde{\lambda} \in \Lambda_H^{(1)}} a(v_{\tilde{\lambda}}, v_{\tilde{\lambda}}) \leq c(k_0) \left\{ a(v_0, v_0) + \sum_{\lambda \in \Lambda} a(v_\lambda, v_\lambda) \right\}.$$

Hence, in view of (29) the lower bound follows from

$$(30) \quad a(v_0, v_0) + \sum_{\lambda \in \Lambda} a(v_\lambda, v_\lambda) \leq C(j+1)^2 a(v, v).$$

Assume for the moment that  $\mathcal{S}_0^\circ \neq \emptyset$ . Then we only have to collect the well-known results of Yserentant [39] to show

$$\begin{aligned} \sum_{\lambda \in \Lambda} a(v_\lambda, v_\lambda) &\leq \alpha_1 \sum_{\lambda \in \Lambda} |v_\lambda|_1^2 \leq c \sum_{k=1}^j 4^k \sum_{\lambda \in \Lambda_k} v(\lambda) \|\lambda\|_0^2 \\ &\leq 2c \sum_{k=1}^j 4^k \|(I_k - I_{k-1})v\|_0^2 \leq C(j+1)^2 |v|_1^2. \end{aligned}$$

In particular, we have employed an inverse inequality ([39, Lem. 3.3]), the boundedness of the incorporated quadrature rule by the  $L^2$ -norm, i.e.,

$$\frac{1}{6} \sum_{t \in \mathcal{T}_k} |t| \sum_{p \in t} |w(p)|^2 \leq 2 \|w\|_0^2, \quad w \in \mathcal{S}_k,$$

and the approximation of unity by the interpolation operators  $I_k$  ([39, Thm. 3.2]). The remaining estimate

$$a(v_0, v_0) = a(I_0 v, I_0 v) \leq C(j+1) |v|_1^2$$

easily follows from the stability of the interpolation ([39, Thm. 3.1]).

As by definition  $\mathcal{S}_0^\circ = \text{span } \hat{\Lambda}_0$ , we still have to consider the case

$$(31) \quad \hat{\Lambda}_{k^*} \neq \emptyset, \quad \hat{\Lambda}_{k^*-1} = \dots = \hat{\Lambda}_0 = \emptyset$$

for some  $k^* > 0$ . Now changing the initial level from 0 to  $k^*$ , we have  $v_\lambda = v(p)\lambda$ ,  $\lambda = \lambda_p^{(k^*)} \in \hat{\Lambda}_{k^*}$ , so that the assertion (30) is immediately obtained from

$$(32) \quad \sum_{p \in \mathcal{N}_{k^*}^\circ} |v(p)\lambda_p^{(k^*)}|_1^2 \leq C |I_{k^*} v|_1^2$$

and the stability of the interpolation cited above. As a consequence of (31), we have  $\mathcal{N}_{k^*-1} = \mathcal{N}_{k^*-1}^\bullet \subset \mathcal{N}_{k^*}^\bullet$  so that (32) follows from Lemma 3.2, with  $\mathcal{N}_{k^*}^\bullet := \mathcal{N}_{k^*}^\bullet$  and  $v_{k^*} := I_{k^*} v \in \mathcal{S}_{k^*}^\bullet := \text{span}\{\lambda \mid \lambda \in \hat{\Lambda}_{k^*}\}$ . This completes the proof of the lower bound of  $a(H_j^{(1)} A_j v, v)$ .

To prove an upper bound by Lemma 3.1(ii), it is sufficient to show that

$$(33) \quad a(v, v) \leq K_1 \left\{ a(\tilde{v}_0, \tilde{v}_0) + \sum_{\tilde{\lambda} \in \Lambda_H^{(1)}} a(v_{\tilde{\lambda}}, v_{\tilde{\lambda}}) \right\}$$

holds for the splitting (28) of some fixed  $v \in \mathcal{S}_j^\circ$ . Recall that the splitting is unique. Using the arguments of the proof of Lemma 3.3, we can show that

$$(34) \quad c(k_0) \sum_{p \in \mathcal{N}_k \cap t} |w(p)|^2 \leq |w|_{1,t}^2 \leq C(k_0) \sum_{p \in \mathcal{N}_k \cap t} |w(p)|^2, \quad t \in \mathcal{T}_k,$$

holds for all  $w \in \text{span } \Lambda_k^{(1)}$ . Based on this norm equivalence, we can extend the proof of the strengthened Cauchy–Schwarz inequality [38, Lem. 2.7] to truncated functions, giving

$$(35) \quad (w_l, w_k)_1 \leq c(k_0) \left( \frac{1}{\sqrt{2}} \right)^{|l-k|-k_0} |w_l|_1 |w_k|_1$$

for all  $w_l \in \text{span } \Lambda_l^{(1)}$ ,  $w_k \in \text{span } \Lambda_k^{(1)}$ , and  $|l-k| \geq k_0$ . From (35) the estimate

$$(36) \quad |v|_1^2 \leq C(k_0) \left\{ |v_0|_1^2 + \sum_{\lambda \in \Lambda_H^{(1)}} |v_\lambda|_1^2 \right\}$$

can be derived by well-known arguments from [38]. Finally, (33) is an immediate consequence of (36) and the ellipticity of  $a(\cdot, \cdot)$ .

By Lemma 3.3 and the equivalence of norms on  $S_0$ , it is obvious that the preconditioner  $\tilde{H}_j^{(1)}$  is just a spectrally equivalent modification of  $H_j^{(1)}$ . This completes the proof of the theorem.  $\square$

For the symmetric preconditioner  $H_j^{(2)}$  we can state a related result without any regularity assumptions imposed on the active set.

**THEOREM 3.5.** *There exist constants  $K_0, K_1$  depending only on  $\alpha_0, \alpha_1$  in (3) and the shape regularity of  $T_0$  such that the estimate*

$$K_0(j + 1)^{-2}a(v, v) \leq a(H_j^{(2)} A_j v, v) \leq K_1 a(v, v)$$

holds for all  $v \in S_j^\circ$ .

*Proof.* Let  $v \in S_j^\circ$ . Based on the unique splitting

$$v = v_0 + \sum_{\lambda \in \Lambda_H^{(2)}} v_\lambda, \quad v_0 \in V_0^{(2)}, \quad v_\lambda \in V_\lambda^{(2)},$$

we can follow the arguments in the proof of Theorem 3.4, with the important difference that the corresponding results on hierarchical bases can be applied directly. Again we have to take care of the case

$$(37) \quad \Lambda_{k^*}^{(2)} \neq \emptyset, \quad \Lambda_{k^*-1}^{(2)} = \dots = \Lambda_0^{(2)} = \emptyset$$

for some  $k^* > 0$ . But as  $\mathcal{N}_{k^*-1}^{\circ, \text{reg}}$  is empty, for each point  $p \in \mathcal{N}_{k^*-1} \cap \mathcal{N}_{k^*}^{\circ, \text{reg}}$  we find at least one  $k^*$ -neighbor  $q \in \mathcal{N}_{k^*} \setminus \mathcal{N}_{k^*-1}$  contained in  $\mathcal{N}_{k^*}^{\bullet, \text{reg}}$ . Hence Lemma 3.2 can be applied as above, setting  $\mathcal{N}_{k^*}^* := \mathcal{N}_{k^*}^{\bullet, \text{reg}}$  and  $S_{k^*}^* := \text{span}\{\lambda \mid \lambda \in \Lambda_{k^*}^{(2)}\}$ .

*Remark 3.4.* Recall that the construction of the preconditioners is independent of the construction of the disjoint splitting  $\mathcal{N}_j = \mathcal{N}_j^\circ \cup \mathcal{N}_j^\bullet$ . In particular, if we are solving an unconstrained elliptic problem, we can define the active set  $\mathcal{N}_j^\bullet$  as the set of all nodes on which the iterative error is considered small enough. A corresponding strategy was proposed in [35]. In this case we cannot expect  $k_0$  in condition (R) to be uniformly bounded (cf. Remark 3.3) so that only the symmetric preconditioner should be used.

*Remark 3.5.* In the proofs of Theorems 3.4 and 3.5 we have extended well-known results on hierarchical bases from the unconstrained to the constrained case by suitable properties of the truncation operators  $T_{j,k}$  or the restriction of the active set  $\mathcal{N}_j^\bullet$ . The same technique can be applied to other multilevel additive Schwarz methods as, for example, in applying the BPX preconditioner to obtain related results in three space dimensions [7].

Theorems 3.4 and 3.5 show that under reasonable assumptions all preconditioners under consideration are spectrally equivalent. However, in the nonsymmetric case the actual constants depend heavily on the constant  $k_0$ , while the behavior of the symmetric preconditioner  $H_j^{(2)}$  has been shown to be more robust with respect to the choice of  $\mathcal{N}_j^\bullet$ . This superiority will be supported by the numerical results presented in §5.

**4. Semi-local and local error estimates.** Let  $u \in H_0^1(\Omega)$  denote the exact solution of (2) and  $u_j \in S_j$  the exact solution of the approximate problem (4) with respect to  $\mathcal{T} = \mathcal{T}_j$ . Expecting that only an approximation  $\tilde{u}_j \in S_j$  of  $u_j$  is known in

actual computations, we are interested in a posteriori error estimates  $\tilde{\varepsilon}$  for the total error  $\varepsilon$ ,

$$\varepsilon = \|u - \tilde{u}_j\| := a(u - \tilde{u}_j, u - \tilde{u}_j)^{1/2},$$

which are efficient and reliable in the sense that

$$(38) \quad \gamma_0 \tilde{\varepsilon} \leq \|u - \tilde{u}_j\| \leq \gamma_1 \tilde{\varepsilon}$$

holds with positive coefficients  $\gamma_0, \gamma_1$  depending only moderately on the refinement level  $j$ . The local contributions to  $\tilde{\varepsilon}$  will be used as local error indicators in the adaptive refinement process. This concept of adaptivity is well established for linear elliptic equations and has been used by a variety of authors. See [3], [14], [23], [28], [37] for further references. Extending the approach of Deuffhard, Leinen, and Yserentant [14], [28] to obstacle problems, we will proceed in two main steps:

*Step 1.* Replace the exact solution  $u$  in (38) by the piecewise quadratic approximation  $U_j \in H_0^1(\Omega)$ .

*Step 2.* Localize the computation of  $U_j$  to obtain  $\tilde{U}_j$  with  $\tilde{\varepsilon} := |\tilde{U}_j - \tilde{u}_j|$  satisfying (38).

The first step is settled by the following lemma, which is a consequence of the triangle inequality.

LEMMA 4.1. *Assume that the piecewise quadratic approximation  $U_j$  is of higher accuracy in the sense that*

$$(39) \quad \|u - U_j\| \leq q \|u - u_j\|, \quad 0 \leq q < 1, \quad j = 0, 1, \dots$$

and  $\tilde{u}_j \in \mathcal{S}_j$  satisfies

$$(40) \quad \|u - u_j\| \leq \sigma \|u - \tilde{u}_j\|, \quad j = 0, 1, \dots$$

with  $q\sigma < 1$  and  $q, \sigma$  not depending on  $j$ . If  $\tilde{\varepsilon}$  satisfies

$$(41) \quad \tilde{\gamma}_0 \tilde{\varepsilon} \leq \|\tilde{u}_j - U_j\| \leq \tilde{\gamma}_1 \tilde{\varepsilon},$$

then (38) holds with  $\gamma_0 = \tilde{\gamma}_0 / (1 + q\sigma)$  and  $\gamma_1 = \tilde{\gamma}_1 / (1 - q\sigma)$ .

*Remark 4.1.* Recall that for sufficiently smooth data the piecewise quadratic approximation is even of higher order than are piecewise linear elements (cf. [12]). In this case (39) is trivial, if the initial triangulation  $\mathcal{T}_0$  is chosen fine enough. Further note that (40) is always satisfied if no obstacle is present because in this case  $u_j$  is the best approximation of  $u$  in  $\mathcal{S}_j$ . In general, (40) follows from

$$\|u_j - \tilde{u}_j\| \leq (1 - 1/\sigma) \|u - u_j\|,$$

with  $\sigma < q^{-1}$ , which may be regarded as an accuracy assumption on  $\tilde{u}_j$ .

In the sequel we assume that (39) and (40) are satisfied to concentrate on the derivation of  $\tilde{\varepsilon}$  with the property (41).

Let  $\mathcal{Q}_j \subset H_0^1(\Omega)$  denote the subspace of piecewise quadratic functions on  $\mathcal{T}_j$  vanishing at the boundary and

$$K_j^{\mathcal{Q}} = \{v \in \mathcal{Q}_j \mid v(p) \leq \varphi^L(p), p \in \mathcal{N}_j, v(e) \leq \varphi^Q(e), e \in \mathcal{E}_j\},$$

the corresponding approximation of the constraints  $K$ . For notation we used  $v(e) := v(\text{midpoint of } e)$ ,  $e \in \mathcal{E}_j$ , for functions  $v : \Omega \rightarrow \mathbb{R}$  and suitable restrictions  $\varphi^L, \varphi^Q$  of the obstacle  $\varphi$  to  $\mathcal{N}_j$  and  $\mathcal{E}_j$ , respectively. Now  $U_j$  can be computed from

$$(42) \quad \text{find } U_j \in K_j^{\mathcal{Q}} \text{ such that } a(U_j, U_j - v) \leq \ell(U_j - v), \quad v \in K_j^{\mathcal{Q}}.$$

For notational convenience the index  $j$  will be suppressed in the following notation. In view of Lemma 4.1 we are interested in the defect  $d = U_j - \tilde{u}_j \in \mathcal{Q}_j$ , which is the unique solution of the following:

$$(43) \quad \text{find } d \in D \text{ such that } a(d, d - v) \leq r(d - v), \quad v \in D.$$

The constraints are given by

$$D = D(\tilde{u}_j) := \{v \in \mathcal{Q}_j \mid v + \tilde{u}_j \in K_j^{\mathcal{Q}}\}$$

and the right-hand side is the residual  $r := \ell - a(\tilde{u}_j, \cdot)$ .

As  $d$  is not available at reasonable computational cost, the remainder of this section will be devoted to the localization of the defect problem (43). A possible way is indicated in the next lemma, showing that (38) is preserved by spectrally equivalent modifications of  $a(\cdot, \cdot)$ .

LEMMA 4.2. *Let  $\tilde{d}$  be the solution of the following:*

$$(44) \quad \text{find } \tilde{d} \in D \text{ such that } \tilde{a}(\tilde{d}, \tilde{d} - v) \leq r(\tilde{d} - v), \quad v \in D$$

with a symmetric form  $\tilde{a}(\cdot, \cdot)$  satisfying

$$(45) \quad c_0 \tilde{a}(v, v) \leq a(v, v) \leq c_1 \tilde{a}(v, v), \quad v \in \mathcal{Q}_j,$$

with positive constants  $c_0, c_1$ . Then

$$(46) \quad C_0 \tilde{a}(\tilde{d}, \tilde{d}) \leq a(d, d) \leq C_1 \tilde{a}(\tilde{d}, \tilde{d})$$

holds with  $C_0 = (c_0^{-1} + 2c_1(1 + c_0^{-1}))^{-1}$ ,  $C_1 = c_1 + 2c_0^{-1}(1 + c_1)$ .

*Proof.* By symmetry arguments it is sufficient to establish the right inequality in (46). Together with (45) we obtain from (43) that

$$a(d, d) \leq c_1 \tilde{a}(\tilde{d}, \tilde{d}) + 2r(d - \tilde{d}).$$

Now the assertion follows from

$$(47) \quad r(d - \tilde{d}) \leq c_0^{-1}(1 + c_1) \tilde{a}(\tilde{d}, \tilde{d}).$$

To show (47), observe that the choice  $v = d$  in (44) leads to

$$(48) \quad r(d - \tilde{d}) \leq \tilde{a}(\tilde{d}, d - \tilde{d}).$$

Hence, in view of Cauchy's inequality it remains to prove

$$(49) \quad |d - \tilde{d}|_{\tilde{a}} \leq c_0^{-1}(1 + c_1) |\tilde{d}|_{\tilde{a}},$$

with  $|\cdot|_{\tilde{a}}$  denoting the energy norm induced by  $\tilde{a}(\cdot, \cdot)$ . It is obvious that  $\tilde{d}$  is the solution of the original problem (43), with  $r$  replaced by a modified right-hand-side  $\tilde{r}$  defined by

$$\tilde{r} := r + a(\tilde{d}, \cdot) - \tilde{a}(\tilde{d}, \cdot).$$

As the solution of variational inequalities depends Lipschitz-continuously on the right-hand side with Lipschitz constant  $c_0^{-1}$  (cf. [25, p. 24]), we obtain (49) from

$$|d - \tilde{d}|_{\tilde{a}} \leq c_0^{-1} \sup_{|v|_{\tilde{a}}=1} |a(\tilde{d}, v) - \tilde{a}(\tilde{d}, v)| \leq c_0^{-1}(1 + c_1) |\tilde{d}|_{\tilde{a}}.$$

This completes the proof.  $\square$

Note that Lemma 4.2 is valid for arbitrary convex constraints and arbitrary space dimensions.

To construct suitable quadratic forms  $\tilde{a}(\cdot, \cdot)$  we introduce the two-level splitting

$$(50) \quad \mathcal{Q}_j = \mathcal{S}^L \oplus \mathcal{S}^Q,$$

which consists of the linear part  $\mathcal{S}^L = \mathcal{S}_j$  and the remaining quadratic part  $\mathcal{S}^Q$ . Note that the quadratic bubbles  $\mu_e \in \mathcal{Q}_j$ ,  $e \in \mathcal{E}_j$ , defined by

$$\mu_e(p) = 0, \quad p \in \mathcal{N}_j, \quad \mu_e(\bar{e}) = \delta_{e, \bar{e}}, \quad \bar{e} \in \mathcal{E}_j,$$

form a basis of  $\mathcal{S}^Q$ . Following (50), we split  $v \in \mathcal{Q}_j$  according to

$$(51) \quad v = v^L + v^Q, \quad v^L \in \mathcal{S}^L, \quad v^Q = \sum_{e \in \mathcal{E}_j} v_e \mu_e \in \mathcal{S}^Q.$$

Then we obtain the quadratic form  $b(\cdot, \cdot)$ ,

$$(52) \quad b(v, w) = a(v^L, w^L) + a^Q(v^Q, w^Q), \quad a^Q(v^Q, w^Q) := \sum_{e \in \mathcal{E}_j} v_e w_e a(\mu_e, \mu_e),$$

by neglecting the coupling of  $\mathcal{S}^L$ ,  $\mathcal{S}^Q$  and  $\mu_e, \mu_g, e \neq g$ , respectively. By also using the preconditioner  $\hat{a}(\cdot, \cdot)$  resulting from the standard hierarchical basis decomposition of  $\mathcal{S}^L = \mathcal{S}_j$ , we end up with

$$(53) \quad \hat{b}(v, w) = \hat{a}(v^L, w^L) + a^Q(v^Q, w^Q).$$

From [14, Lem., p. 14] and the following considerations it is well known that

$$(54) \quad cb(v, v) \leq a(v, v) \leq C(j + 1)^2 \hat{b}(v, v), \quad v \in \mathcal{Q}_j,$$

holds with suitable constants  $c, C$ . Summarizing these results, we obtain the first important result of this section.

**THEOREM 4.3.** *Assume that the conditions (39) and (40) are satisfied. Let  $\hat{d}$  be the solution of the semilocal problem*

$$(55) \quad \text{find } \hat{d} \in D \text{ such that } \hat{b}(\hat{d}, \hat{d} - v) \leq r(\hat{d} - v), \quad v \in D.$$

Then (38) holds for

$$\tilde{\varepsilon}^2 := |\hat{d}|_b^2 = \hat{a}(\hat{d}^L, \hat{d}^L) + a^Q(\hat{d}^Q, \hat{d}^Q)$$

and constants  $\gamma_0 = \hat{\gamma}_0/(j + 1)$ ,  $\gamma_1 = \hat{\gamma}_1(j + 1)$ . Here  $\hat{\gamma}_0, \hat{\gamma}_1$  are depending only on  $q\sigma$ , the ellipticity of  $a(\cdot, \cdot)$ , and the shape regularity of  $\mathcal{T}_0$ .

*Proof.* Using (54), Theorem 4.3 is an immediate consequence of Lemmas 4.1 and 4.2.  $\square$

**Remark 4.2.** The error estimate (55) is called *semilocal* because the frequencies of  $\hat{d}$  are decoupled with respect to the quadratic form but coupled by the set of constraints  $D$ . In our numerical experiments we will use the local contributions

$$\eta_e = (\hat{d}_e^Q)^2 a(\mu_e, \mu_e), \quad e \in \mathcal{E}_j,$$

of  $a^Q(\hat{d}^Q, \hat{d}^Q)$  as local error indicators in the adaptive refinement process. Of course, (55) reduces to the error estimate proposed in [14], if the obstacle is not active.

*Remark 4.3.* The simplified defect problem (55) may be solved approximately using the active-set strategy described above. Because the preconditioners proposed in the preceding section are just truncated versions of  $\hat{a}(\cdot, \cdot)$ , we can expect the corresponding linear subproblems to be solved very efficiently.

To derive a less robust but local error estimate, we consider the simplified defect problem:

$$(56) \quad \text{find } \delta \in D \text{ such that } b(\delta, \delta - v) \leq r(\delta - v), \quad v \in D.$$

Recall that

$$(57) \quad c_0 b(v, v) \leq a(v, v) \leq c_1 b(v, v), \quad v \in \mathcal{Q}_j,$$

with positive constants  $c_0, c_1$  independent of  $j$  (cf. [14]). Assuming (39) and (40), it follows from Lemmas 4.1 and 4.2 that the solution  $\delta$  of (56) provides an error estimate with the property (38). Now (56) is decoupled by one block Gauss–Seidel iteration step applied to the initial iterate zero, i.e., we compute an estimate  $\hat{\delta} = \hat{\delta}^L + \hat{\delta}^Q$  from

$$(58) \quad \text{find } \hat{\delta}^L \in D^L \text{ such that } a(\hat{\delta}^L, \hat{\delta}^L - v) \leq r^L(\hat{\delta}^L - v), \quad v \in D^L,$$

and

$$(59) \quad \begin{aligned} &\text{find } \hat{\delta}^Q \in D^Q(\hat{\delta}^L) \text{ such that} \\ &a^Q(\hat{\delta}^Q, \hat{\delta}^Q - v) \leq r^Q(\hat{\delta}^Q - v), \quad v \in D^Q(\hat{\delta}^L), \end{aligned}$$

where  $r^L, r^Q$  denote the restriction of  $r$  to  $\mathcal{S}^L, \mathcal{S}^Q$  and  $D^L, D^Q(\hat{\delta}^L)$  are defined by

$$D^L = \mathcal{S}^L \cap D, \quad D^Q(w^L) = \{v^Q \in \mathcal{S}^Q \mid v^Q + w^L \in D\}, \quad w^L \in \mathcal{S}^L.$$

Assuming that

$$K_j = \{v \in \mathcal{S}_j \mid v(p) \leq \varphi^L(p), p \in \mathcal{N}_j\} \subset K_j^Q,$$

the linear defect problem is recovered by (58) with the consequence that

$$\hat{\delta}^L = u_j - \tilde{u}_j.$$

Moreover, each component  $\hat{\delta}_e^Q$  of  $\hat{\delta}^Q$  can be computed separately, giving

$$(60) \quad \hat{\delta}_e^Q = \min\{r^Q(\mu_e)/a(\mu_e, \mu_e), (\varphi^Q - \hat{\delta}^L - \tilde{u}_j)(e)\}, \quad e \in \mathcal{E}_j.$$

Hence

$$(61) \quad \tilde{\varepsilon}^2 := |\hat{\delta}|_b^2 = \|u_j - \tilde{u}_j\|^2 + a^Q(\hat{\delta}^Q, \hat{\delta}^Q)$$

provides a local error estimate as soon as the iterative error  $\|u_j - \tilde{u}_j\|$  is known. Again (61) reduces to the error estimate proposed in [14] if the obstacle is not active, and the local contributions to  $a^Q(\hat{\delta}^Q, \hat{\delta}^Q)$  may be used as local error indicators in the adaptive refinement process.

We will make use of the interpolation operator  $\pi : \mathcal{S}^L \rightarrow \mathcal{S}^Q$ , defined by

$$\pi(v^L)(p_e) = (v^L(p_1) + v^L(p_2))/2, \quad e = (p_1, p_2) \in \mathcal{E}_j, \quad v^L \in \mathcal{S}^L,$$

to show that (61) provides a lower bound for the total error.

**THEOREM 4.4.** *Assume that the conditions (39) and (40) are satisfied. Let  $K_j \subset K_j^Q$  and assume that*

$$(L) \quad |\pi(\delta^L - \hat{\delta}^L)|_{aQ} \leq \beta \|\delta^L - \hat{\delta}^L\|$$

holds with a positive constant  $\beta$  independent of  $j$ . Then

$$(62) \quad \gamma_0 |\hat{\delta}|_b \leq \|\tilde{u}_j - u\|$$

holds with a positive constant  $\gamma_0$  depending only on  $q\sigma$ ,  $\beta$ , the ellipticity of  $a(\cdot, \cdot)$ , and the shape regularity of  $\mathcal{T}_0$ .

*Proof.* First recall that we have from Lemmas 4.1 and 4.2

$$(63) \quad c_0 |\delta|_b \leq \|u - \tilde{u}_j\| \leq c_1 |\delta|_b,$$

with constants  $c_0, c_1$  independent of  $j$ . As  $K_j \subset K_j^Q$ , we obtain  $\hat{\delta}^L = u_j - \tilde{u}_j$  so that (40) leads to

$$(64) \quad \|\hat{\delta}^L\| \leq (\sigma + 1) \|u - \tilde{u}_j\| \leq c |\delta|_b.$$

The estimation of the quadratic part of  $|\hat{\delta}|_b^2 = \|\hat{\delta}^L\|^2 + |\hat{\delta}^Q|_{aQ}^2$  is more complicated. Obviously  $\delta^Q$  is the solution of

$$(65) \quad \text{find } \delta^Q \in D^Q(\delta^L) \text{ such that } a^Q(\delta^Q, \delta^Q - v) \leq r(\delta^Q - v), \quad v \in D^Q(\delta^L),$$

with  $\delta = \delta^L + \delta^Q$ . By representing (65) as a complementary problem it is easily verified that (59) and (65) are symmetric with respect to the obstacle and the right-hand side. More precisely, (59) and (65) can be replaced by

$$(66) \quad \text{find } \hat{\delta}^Q \in R \text{ such that } a^Q(\hat{\delta}^Q, \hat{\delta}^Q - v) \leq a^Q(\varphi^Q - \pi(\tilde{u}_j + \hat{\delta}^L), \hat{\delta}^Q - v), \quad v \in R,$$

and

$$(67) \quad \text{find } \delta^Q \in R \text{ such that } a^Q(\delta^Q, \delta^Q - v) \leq a^Q(\varphi^Q - \pi(\tilde{u}_j + \delta^L), \delta^Q - v), \quad v \in R,$$

with constraints

$$R = \{v \in \mathcal{S}^Q \mid v_e \leq r^Q(\mu_e)/a(\mu_e, \mu_e), e \in \mathcal{E}_j\}.$$

Again we assume that (66), (67) are Lipschitz with respect to the right-hand side in order to obtain by assumption (L) the following inequality:

$$(68) \quad |\delta^Q - \hat{\delta}^Q|_{aQ} \leq |\pi(\delta^L - \hat{\delta}^L)|_{aQ} \leq \beta \|\delta^L - \hat{\delta}^L\|.$$

Now the triangle inequality gives

$$(69) \quad |\hat{\delta}^Q|_{aQ}^2 \leq 4\beta^2 (\|\hat{\delta}^L\|^2 + |\delta|_b^2)$$

and the assertion follows, along with (63) and (64).  $\square$

*Remark 4.4.* In the simplified case of a quasiuniform sequence of triangulations with meshsize  $h_j$  the condition (L) is equivalent to

$$(70) \quad \|\delta^L - \hat{\delta}^L\|_0 \leq ch_j |\delta^L - \hat{\delta}^L|_1,$$

with  $c$  independent of  $j$ . Obviously (70) is always satisfied with  $c = c(j)$ , giving (62) with  $\gamma_0 = \gamma_0(j)$ . In general, (70) may be regarded as a regularity condition on  $\delta$ . Indeed, assuming  $\tilde{u}_j = u_j$  and regarding  $\delta^L$  as a perturbation of  $\hat{\delta}^L = 0$  by the coupling with  $\delta^Q$  at the free boundary, condition (L) is satisfied if these perturbations remain local with increasing  $j$ .

Of course, a further restriction is imposed by the assumption  $K_j \subset K_j^Q$ , which, for example, is satisfied if the obstacle function is continuous and piecewise linear and the initial triangulation is chosen appropriately (cf., e.g., the example treated in the following section).

The error estimate (61) was originally proposed in [26] and [27] for the adaptive solution of a special obstacle problem arising in semiconductor device simulation. In this special problem we can expect from the physical data that the error is dominated uniformly in  $j$  by contributions generated away from the free boundary, suffering only minor effects from the localization (58), (59). In particular, the nonactive region can always be resolved with sufficient accuracy on the initial triangulation  $\mathcal{T}_0$ . Under these assumptions we can easily prove that (61) is reliable in the sense of (38), particularly in that it also provides a uniform upper bound of the exact error  $\varepsilon$ .

However, simple examples show that (61) may deliver  $\hat{\delta} = 0$  even though  $d \neq 0$  holds true. Together with Theorem 4.4 this indicates that (61) is likely to underestimate the true error, which will be confirmed by numerical experiments reported in the next section.

**5. Numerical results.** In this section we concentrate on composing an adaptive Multilevel Method from the modules described above. This method is then applied to a challenging model problem confirming the properties expected from the theoretical considerations.

On each refinement level  $j$  we apply the active-set strategy given in §2 until the active set is left invariant. The iteration is started with the interpolated approximation from the previous level, with the value at each node having at least one active neighbor projected to the obstacle. On the first level the obstacle function is used as the initial iterate. Each step of the outer iteration requires the solution of the linear subproblem (11), which is performed iteratively by cg iterations preconditioned by the reduced hierarchical basis preconditioners introduced above. This inner iteration is stopped as soon as the estimated linear iteration error  $\kappa$  satisfies  $\kappa \leq \kappa_0$ . Here estimate  $\kappa$  is computed as described in [14]. Recall that the threshold  $\kappa_0$  has to be chosen small enough to ensure the convergence of the outer iteration (cf. Remark 2.1). In the following example,  $\kappa_0 = 10^{-3}$  is used.

The same algorithm with  $\kappa_0$  replaced by  $\kappa'_0 = 10^{-2}$  is applied to the solution of the semilocal defect problem (55), providing the error estimate  $\varepsilon^s = |\hat{d}|_b$ . A local error estimate  $\varepsilon^l = |\hat{\delta}|_b$  is obtained by using the iterative error of the final linear subproblem as an approximation for  $\hat{\delta}^L = u_j - \tilde{u}_j$  and evaluating (59). The iterative solution of the semilocal defect problem is started with the local estimate  $(0, \hat{\delta}^Q)$ .

According to Remark 4.2, an edge  $e \in \mathcal{E}_j$  is bisected if its contribution  $\eta_e$  exceeds a certain threshold  $\bar{\eta}$ . To determine  $\bar{\eta}$  we extrapolate  $\eta_e$  as proposed in [2] (see [26]

for details). A new triangulation is constructed by red refinements and green closures (refer to [3], [28], [33], and [34] for further information).

Now we apply the algorithm to a well-known problem describing the elastoplastic torsion of a cylindrical bar with quadratical cross section  $\Omega = (0, 1) \times (0, 1)$ , which is twisted at its upper end around the longitudinal axis in such a way that the lateral surface remains stress free. By modelling the plastic region according to the von Mises yield criterion and normalizing physical constants, it has been shown in [11] that for positive twist angle  $C$  per unit length the stress potential  $u$  is the solution of the variational inequality (2) with  $a(\cdot, \cdot)$ ,  $\ell(\cdot)$  given by

$$a(v, w) = \int_{\Omega} (\partial_1 v \partial_1 w + \partial_2 v \partial_2 w) dx, \quad \ell(v) = 2C \int_{\Omega} v dx$$

and constraints  $K$ ,

$$K = \{v \in H_0^1(\Omega) \mid v(x) \leq \text{dist}(x, \partial\Omega), \text{ a.e. in } \Omega\}.$$

The active points characterize the plastic region while the material is considered elastic in nonactive points. We refer to [17] and [19] for the numerical treatment and to [32] for a theoretical analysis of the problem.

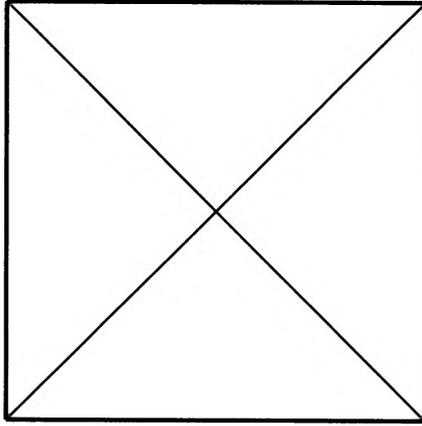
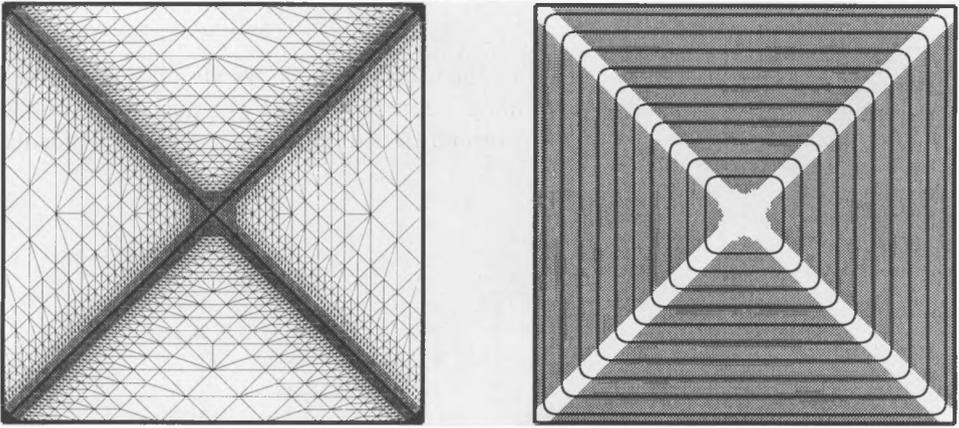
Note that the problem has singular perturbation character with respect to the elastic region, which is located along the diagonals and shrinks for increasing  $C$ .

TABLE 1  
*Iteration history.*

Level	Depth	Nodes	Iterations	
			Solution	Error Estimate
0	0	5	1/0.0	2/0.5
1	1	13	1/0.0	2/1.0
2	2	29	1/0.0	3/1.0
3	3	57	2/0.5	3/2.3
4	4	153	2/2.5	3/3.6
5	5	381	2/5.0	4/2.0
6	5	541	3/3.0	3/2.0
7	5	749	3/3.3	1/0.0
8	6	1605	3/4.3	2/0.0
9	7	5793	4/5.5	2/0.0
10	8	6265	3/6.0	2/0.0

Starting with the initial triangulation  $\mathcal{T}_0$  depicted in Fig. 2 and choosing  $C = 15$ , all nodal points remain active up to the third (uniform) refinement level, rendering a quite challenging problem for an adaptive multilevel method.

In Table 1 we report the number of iterations required by the solution process. The data are presented in the form “number of outer iterations/average number of inner iterations” both needed for the solution and the semilocal error estimate, respectively. In both cases the symmetric version of the hierarchical basis preconditioner is used. The difficulty of detecting the elastic region leads to the difference between depth and refinement level arising from level 5 to level 7. In the sequel the actual number of refinement levels is indicated by subscript in spite of some ambiguity compared to the notation in §3. Note that  $\mathcal{T}_7$  finally allows for a satisfying resolution of the elastic zone. Up to this level the computational work is dominated by the error estimation, providing the local error indicators for the adaptive refinement process. On the subsequent levels the semilocal error estimate automatically reduces to the

FIG. 2. *Initial triangulation  $\mathcal{T}_0$ .*FIG. 3. *Final triangulation  $\mathcal{T}_{10}$  and solution  $U_{10}$ .*

local error estimate. Indeed, the outer iterations do not change the initial guess and may be skipped.

The final triangulation  $\mathcal{T}_{10}$  is depicted in Fig. 3 along with the level curves and the elastic region of the corresponding solution.

The behavior of both error estimates is illustrated in more detail in Fig. 4. Again it is obvious that the situation changes at level 7 (749 nodes), showing a significant decrease of the “exact” error, and both estimates. To compute the “exact” error, we performed a uniform refinement of  $\mathcal{T}_{10}$  and computed the difference to the corresponding solution. Note that only the semilocal estimate provides satisfactory results on lower levels. In fact, due to the very coarse initial grid the local error estimate fails in this example, providing  $\varepsilon_j^l = 0$  for  $j = 0, 1, 2$ . Recall that the performance of both error estimates could be expected from the theoretical considerations in the preceding section. In particular, the local estimate (59) should not be used until the underlying triangulation is fine enough to detect all parts of the inactive region but works very effectively from this moment on.

The final Fig. 5 gives a comparison of both versions of the hierarchical basis preconditioners. To amplify the different behavior we choose  $\kappa_0$  very small, i.e.,

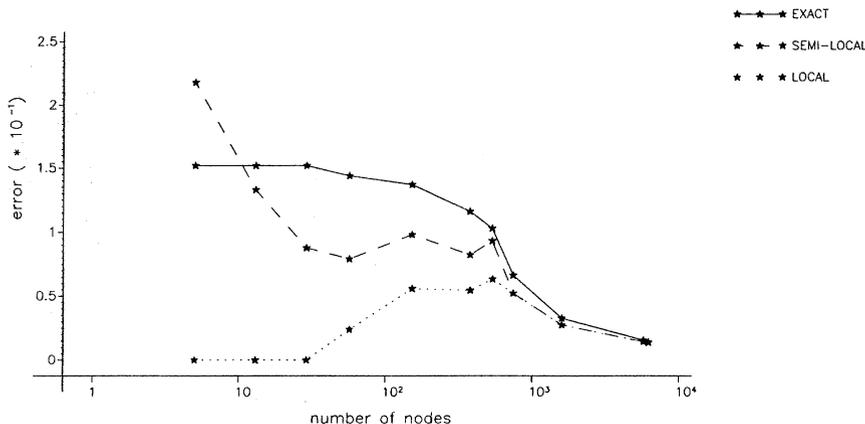


FIG. 4. Comparison of the error estimates.

$\kappa_0 = 10^{-8}$  and the initial iterate is fixed to the upper obstacle for all inner iterations. For each refinement level we choose the linear subproblem with the maximal number of unknowns and report the number of (preconditioned) cg iterations required for its solution.

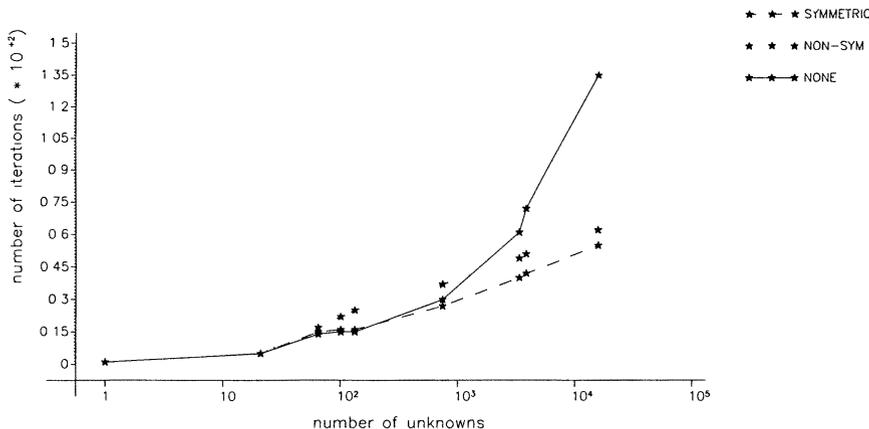


FIG. 5. Comparison of the preconditioners.

As expected, multilevel preconditioning does not improve the convergence of the cg iteration as long as the actual problem allows no suitable representation on the coarser triangulations. Obviously, the nonsymmetric version even causes deterioration of the convergence until the contribution of nontruncated hierarchical basis functions becomes dominant on level 9. On the other hand, the symmetric version immediately takes advantage of the good resolution on level 7 (133 unknowns) and does not lead to deterioration of the convergence on lower levels. Note that in both cases the number of iterations becomes a linear function of the refinement level  $j$ , if  $j$  is large enough. This is exactly the behaviour predicted by the theoretical results derived in §3.

**Acknowledgment.** The authors wish to thank F. Bornemann, R. Roitzsch, H. Yserentant, and the referees for various important remarks and suggestions and S. Wacker for her careful  $\text{\TeX}$ -typing of the manuscript. Special thanks to the production editor J. Anderson from SINUM for her patience and diligence in the course of the publication.

## REFERENCES

- [1] O. AXELSSON AND V. A. BARKER, *Finite Element Solution of Boundary Value Problems*, Academic Press, New York, 1984.
- [2] I. BABUŠKA AND W. C. RHEINBOLDT, *Estimates for adaptive finite element computations*, *SIAM J. Numer. Anal.*, 15 (1978), pp. 736–754.
- [3] R. E. BANK, *PLTMG—A Software Package for Solving Elliptic Partial Differential Equations. User's Guide 6.0.*, Society for Industrial and Applied Mathematics, Philadelphia, PA, 1990.
- [4] R. E. BANK, T. DUPONT, AND H. YSERENTANT, *The hierarchical basis multigrid method*, *Numer. Math.*, 52 (1988), pp. 387–404.
- [5] R. E. BANK, A. H. SHERMAN, AND H. WEISER, *Refinement algorithms and data structures for regular local mesh refinement*, in *Scientific Computing*, R. Stepleman et al., eds., IMACS North-Holland, Amsterdam, 1983, pp. 3–17.
- [6] F. BORNEMANN, *A sharpened condition number estimate for the BPX preconditioner of elliptic finite element problems on highly nonuniform triangulations*, preprint SC 91–9, Konrad-Zuse-Zentrum Berlin, Berlin, 1991.
- [7] F. BORNEMANN, B. ERDMANN, AND R. KORNHUBER, *Adaptive multilevel-methods in three space dimensions*, *Internat. J. Numer. Methods Engrg.*, 36 (1993), pp. 3187–3203.
- [8] F. BORNEMANN AND H. YSERENTANT, *A basic norm equivalence for the theory of multilevel methods*, *Numer. Math.*, 64 (1993), pp. 445–476.
- [9] J. H. BRAMBLE, J. E. PASCIAK, AND J. XU, *Parallel multilevel preconditioners*, *Math. Comp.*, 55 (1990), pp. 1–22.
- [10] A. BRANDT AND C.W. CRYER, *Multigrid algorithms for the solution of linear complementarity problems arising from free boundary problems*, *SIAM J. Sci. Statist. Comput.*, 4 (1983), pp. 655–684.
- [11] H. BRÉZIS AND M. SIBONY, *Equivalence de deux inéquations variationnelles et applications*, *Arch. Rational Mech. Anal.*, 41 (1971), pp. 254–265.
- [12] F. BREZZI, W. W. HAGER, AND P. A. RAVIART, *Error estimates for the finite element solution of variational inequalities I*, *Numer. Math.*, 28 (1977), pp. 431–443.
- [13] R. W. COTTLE, J. S PANG, AND R. E. STONE, *The Linear Complementary Problem*, Academic Press, New York, 1992.
- [14] P. DEUFLHARD, P. LEINEN, AND H. YSERENTANT, *Concepts of an adaptive hierarchical finite element code*, *IMPACT Comput. Sci. Engrg.*, 1 (1989), pp. 3–35.
- [15] M. DRYJA AND O. B. WIDLUND, *Multilevel additive methods for elliptic finite element problems*, *Tech. Report 507*, Courant Institute of Mathematical Sciences, New York, 1990.
- [16] ———, *Towards a unified theory of domain decomposition algorithms for elliptic problems*, in *Domain Decomposition Methods for Partial Differential Equations*, T.F. Chan et al., eds., Society for Industrial and Applied Mathematics, Philadelphia, 1989, pp. 3–21.
- [17] R. GLOWINSKI, J. L. LIONS, AND R. TRÉMOLIÈRES, *Numerical Analysis of Variational Inequalities*, North-Holland, Amsterdam, 1981.
- [18] W. HACKBUSCH AND H. D. MITTELMANN, *On multigrid methods for variational inequalities*, *Numer. Math.*, 42 (1983), pp. 65–76.
- [19] R. H. W. HOPPE, *Multigrid algorithms for variational inequalities*, *SIAM J. Numer. Anal.*, 24 (1987), pp. 1046–1065.
- [20] ———, *Two-sided approximations for unilateral variational inequalities by multigrid methods*, *Optimization*, 18 (1987), pp. 867–881.
- [21] ———, *Une méthode multigrille pour la solution des problèmes d'obstacle*, *M<sup>2</sup> Math. Model. Numer. Anal.*, 24 (1990), pp. 711–736.
- [22] R. H. W. HOPPE AND R. KORNHUBER, *Multilevel preconditioned cg iterations for variational inequalities*, in *Preliminary Proceedings of the 5th Copper Mountain Conference on Multigrid Methods*, P. Frederickson, J. Mandel, and S. McCormick, eds., Copper Mountain Colorado, 1991.
- [23] C. JOHNSON, *Numerical Solutions of Partial Differential Equations by the Finite Element Method*, Cambridge University Press, Cambridge, 1987.

- [24] C. JOHNSON, *Adaptive finite element methods for the obstacle problem*, preprint NO 1991-25, Chalmers University of Technology, Göteborg, Sweden, 1991.
- [25] D. KINDERLEHRER AND G. STAMPACCHIA, *An Introduction to Variational Inequalities and Their Applications*, Academic Press, New York, 1980.
- [26] R. KORNUBER AND R. ROITZSCH, *Self adaptive finite element simulation of bipolar, strongly reverse biased pn-junctions*, Comm. Numer. Methods Engrg., 9 (1993), pp. 243-250.
- [27] ———, *Self adaptive computation of the breakdown voltage of planar pn-junctions with multistep field plates*, in Proceedings of the 4th International Conference of Simulation of Semiconductor Devices and Processes, W. Fichtner and D. Aemmer, eds., Zurich, Switzerland, 1991, pp. 535-543.
- [28] P. LEINEN, *Ein schneller adaptiver Löser für elliptische Randwertprobleme*, Ph.D. thesis, Dortmund, Germany, 1990.
- [29] P. L. LIONS, *On the Schwarz alternating method I*, in Domain Decomposition Methods for Partial Differential Equations, R. Glowinski et al., eds., Society for Industrial and Applied Mathematics, Philadelphia, PA, 1988.
- [30] J. MANDEL, *Etude algébrique d'une méthode multigrille pour quelques problèmes de frontière libre*, C.R. Acad. Sci. Paris, Ser. I, 298 (1984), pp. 469-472.
- [31] ———, *A multilevel iterative method for symmetric, positive definite linear complementarity problems*, Appl. Math. Optimization, 11 (1984), pp. 77-95.
- [32] J.-F. RODRIGUES, *Obstacle Problems in Mathematical Physics*, North-Holland Mathematical Studies 134, North-Holland, Amsterdam, 1987.
- [33] R. ROITZSCH, *KASKADE User's Manual*, Tech. Report TR 89-4, Konrad-Zuse-Zentrum Berlin, Berlin, Germany, 1989.
- [34] ———, *KASKADE Programmer's Manual*, Tech. Report TR 89-5, Konrad-Zuse-Zentrum Berlin, Berlin, Germany, 1989.
- [35] U. RÜDE, *Fully adaptive multigrid methods*, SIAM J. Numer. Anal., 30 (1993), pp. 230-248.
- [36] R. SCHWENKERT *Lösung linearer Komplementaritätsprobleme mittels hierarchischer Basen*, Report 108 Institut für Angewandte Mathematik und Statistik der Universität Würzburg, Würzburg, Germany, 1988.
- [37] B. SZABÓ AND I. BABUŠKA, *Finite Element Analysis*, Wiley & Sons, New York, 1991.
- [38] H. YSERENTANT, *On the multilevel splitting of finite element spaces*, Numer. Math., 49 (1986), pp. 379-412.
- [39] ———, *Two preconditioners based on the multilevel splitting of finite element spaces*, Numer. Math., 58 (1990), pp. 163-184.
- [40] ———, Private communication, 1990.
- [41] X. ZHANG, *Multilevel Schwarz methods*, Numer. Math., 63 (1992), pp. 521-539.