

## The Construction of Preconditioners for Elliptic Problems by Substructuring, III\*

By James H. Bramble, Joseph E. Pasciak, and Alfred H. Schatz

**Abstract.** In earlier parts of this series of papers, we constructed preconditioners for the discrete systems of equations arising from the numerical approximation of elliptic boundary value problems. The resulting algorithms are well suited for implementation on computers with parallel architecture. In this paper, we will develop a technique which utilizes these earlier methods to derive even more efficient preconditioners. The iterative algorithms using these new preconditioners converge to the solution of the discrete equations with a rate that is independent of the number of unknowns. These preconditioners involve an incomplete Chebyshev iteration for boundary interface conditions which results in a negligible increase in the amount of computational work. Theoretical estimates and the results of numerical experiments are given which demonstrate the effectiveness of the methods.

**1. Introduction.** The aim of this series of papers is to propose and analyze methods for efficiently solving the equations resulting from finite element discretizations of second-order elliptic boundary value problems on general domains in  $R^2$  and  $R^3$ . In particular, we shall be concerned with constructing easily invertible and “effective” preconditioners for the resulting system of discrete equations which can be used in a preconditioned iterative algorithm to define a rapid solution method. The methods developed are well suited to parallel computing architectures.

In Parts I and II (references [4] and [5]), we described and analyzed methods for constructing preconditioners for elliptic boundary value problems on polygonal domains in  $R^2$  and  $R^3$ . The proposed methods were based on decomposing the domain into subdomains of size  $d$  and involved the solution of related problems on the subdomains and lower-order coupling systems on the subdomain boundaries. The condition number for the preconditioned system was shown to be on the order of  $(1 + \ln(d/h))^2$  for the method of [4] and  $d/h$  for the method of [5]. Here  $h$  is the mesh size. In this paper, we describe a technique which can utilize such methods to develop more efficient preconditioners. The condition numbers for the resulting preconditioned systems will be made independent of  $d$  and  $h$  with only a slight increase in computational effort.

---

Received May 7, 1987.

1980 *Mathematics Subject Classification* (1985 *Revision*). Primary 65N30; Secondary 65F10.

\*This manuscript has been authored under contract number DE-AC02-76CH00016 with the U.S. Department of Energy. Accordingly, the U.S. Government retains a non-exclusive, royalty-free license to publish or reproduce the published form of this contribution, or allow others to do so, for U.S. Government purposes. This work was also supported in part under the National Science Foundation Grant No. DMS84-05352 and under the Air Force Office of Scientific Research, Contract No. ISSA86-0026 and by the U.S. Army Research Office through the Mathematical Science Institute, Cornell University.

Let  $\Omega$  be a bounded domain in  $R^2$  with a piecewise smooth boundary  $\partial\Omega$ . As a model problem for a second-order uniformly elliptic equation, we shall consider the Dirichlet problem

$$(1.1) \quad \begin{aligned} Lu &= f \quad \text{in } \Omega, \\ u &= 0 \quad \text{on } \partial\Omega, \end{aligned}$$

where

$$Lv = - \sum_{i,j=1}^2 \frac{\partial}{\partial x_i} \left( a_{ij} \frac{\partial v}{\partial x_j} \right),$$

with  $a_{ij}$  symmetric and uniformly positive definite, bounded and piecewise smooth on  $\Omega$ . The generalized Dirichlet form is given by

$$(1.2) \quad A(v, \phi) = \sum_{i,j=1}^2 \int_{\Omega} a_{ij} \frac{\partial v}{\partial x_i} \frac{\partial \phi}{\partial x_j} dx,$$

which is defined for all  $v$  and  $\phi$  in the Sobolev space  $H^1(\Omega)$  (the space of distributions with square-integrable first derivatives). The  $L^2(\Omega)$  inner product is denoted

$$(v, \phi) = \int_{\Omega} v \phi dx.$$

The subspace  $H_0^1(\Omega)$  is the completion of the smooth functions with support in  $\Omega$  with respect to the norm in  $H^1(\Omega)$ . The weak formulation of the problem defined by (1.1) is: Find  $u \in H_0^1(\Omega)$  such that

$$(1.3) \quad A(u, \phi) = (f, \phi) \quad \text{for all } \phi \in H_0^1(\Omega).$$

This leads immediately to the standard Galerkin approximation. Let  $S_h^0(\Omega)$  be a finite-dimensional subspace of  $H_0^1(\Omega)$ . The Galerkin approximation is defined as the solution of the following problem: Find  $U \in S_h^0(\Omega)$  such that

$$(1.4) \quad A(U, \Phi) = (f, \Phi) \quad \text{for all } \Phi \in S_h^0(\Omega).$$

The underlying method which we will consider is a preconditioned iterative method. As explained in Part I, the task of defining a preconditioner for the matrix problem corresponding to (1.4) is the same as that of defining another positive definite form  $B(\cdot, \cdot)$  on  $S_h^0(\Omega) \times S_h^0(\Omega)$ . The importance of making a “good” choice for  $B$  is well known. The form  $B$  will define a good preconditioner provided it has two basic properties. First, the problem of finding the function  $W \in S_h^0(\Omega)$  satisfying

$$(1.5) \quad B(W, \Phi) = G(\Phi) \quad \text{for all } \Phi \in S_h^0(\Omega),$$

for a given linear functional  $G$ , should be more economical to solve on a given computer architecture than (1.4). Secondly,  $B$  should be spectrally close to  $A$  in the sense that there are positive numbers  $\beta_0$  and  $\beta_1$  satisfying

$$(1.6) \quad \beta_0 B(V, V) \leq A(V, V) \leq \beta_1 B(V, V) \quad \text{for all } V \in S_h^0(\Omega),$$

where the ratio  $\beta_1/\beta_0$  is not too large. These two properties will guarantee, firstly, that the work per iterative step in applying the preconditioned method will be small, and, secondly, that the number of steps to reduce the error to a given size will also be small, so that an efficient algorithm will result.

In Section 2 the form  $B$  is defined and the essential step in the iterative algorithm of solving (1.5) is described. This form is defined in terms of a polynomial  $P_m$  which is related to the classical Chebyshev polynomials. The relevant properties are given in Section 3. Section 4 discusses various computational aspects of the method in a more general setting. Finally, in Section 5, the results of numerical experiments are given. These computations show that the theoretical estimates are fully realized in practice.

For other works dealing with the numerical solution of boundary value problems via substructuring we refer the reader to [1], [2], [3], [6], [7], [8], [9]. We emphasize that a novel feature of our methods [4], [5] is that more than two subdomains can meet at an interior point of the original domain. For example, our methods apply to a checkerboard subdivision of a square. Using the technique of this paper, the condition number for the resulting system is shown to be bounded independently of the number of such points.

**2. The Construction of  $B(\cdot, \cdot)$  and the Preconditioning Algorithm.** As mentioned in the introduction, the preconditioner which we will construct involves the solution of smaller related problems on subdomains and subdomain boundaries. As in Part I, for the sake of simplicity of exposition, we shall proceed with the discussion only for the special case of polygonal domains and piecewise linear approximations.

More precisely, we shall begin with the following assumptions with regard to  $\Omega$ . These assumptions are the same as those given in Section 2 of Part I, and hence the results given in Section 3 of Part I apply.

A.1:  $\Omega$  is a polygonal domain.

A.2: For each  $h$ ,  $0 < h < 1$  a parameter,  $\Omega$  has been given a quasi-uniform triangulation  $\Omega^h$ . By this we mean that there exists a positive constant  $c_1$  independent of  $h$  such that each triangle  $\tau^h \in \Omega^h$  contains a ball of radius  $c_1 h$  and is contained in a ball of radius  $h$ .

A.3: For each triangulation  $\Omega^h$ ,  $\Omega$  may be written in terms of  $n_r$  disjoint regions,  $\Omega_k$ , with  $\bar{\Omega} = \bigcup \bar{\Omega}_k$ , which are either quadrilaterals or triangles whose sides coincide with the mesh lines of the original triangulation and which are quasi-uniform of size  $d \geq h$  with constants as above which are independent of  $d$  and  $h$ . If  $\Omega_k$  is a quadrilateral, we require additionally that the lengths of each side be bounded from below by  $c_1 d$  and that any interior angle  $\alpha$  satisfy  $0 < C_0 \leq \alpha \leq C_1 < \pi$ . The collection of regions  $\Omega_k$  will frequently be referred to as the subdomains (see Figure 2.1).

For each  $h$ , let  $S_h(\Omega)$  be the space of continuous piecewise linear functions defined relative to the triangulation  $\Omega^h$  and  $S_h^0(\Omega)$  be the subspace of  $S_h(\Omega)$  consisting of those functions which vanish on  $\partial\Omega$ .  $S_h^0(\Omega_k)$  will denote the subspace of  $S_h^0(\Omega)$  of functions whose supports are contained in  $\bar{\Omega}_k$  (in particular, they vanish on  $\partial\Omega_k$  and outside  $\bar{\Omega}_k$ ). In addition, let  $S_h(\Omega_k)$  be the set of functions which are restrictions of those in  $S_h^0(\Omega)$  to  $\bar{\Omega}_k$ .  $S_h(\partial\Omega_k)$  will denote the restrictions of  $S_h(\Omega_k)$  to  $\partial\Omega_k$ . Let  $\Gamma = \bigcup \partial\Omega_j$  and  $S_h(\Gamma)$  be the restriction of functions in  $S_h^0(\Omega)$  to  $\Gamma$ . In what follows,  $c$  and  $C$  (with or without subscript) will denote generic positive constants which are independent of  $h$ ,  $d$  and the regions  $\Omega_k$ .

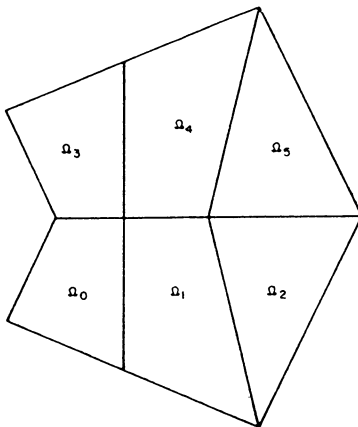


FIGURE 2.1

*A typical domain with subdomains.*

For simplicity of presentation, we shall restrict our development to the case in which each  $\partial\Omega_k$  has a uniformly (with respect to arc length) spaced grid. This restriction will be removed in Section 4. We define

$$A_k(u, v) = \sum_{i,j=1}^2 \int_{\Omega_k} a_{ij} \frac{\partial u}{\partial x_i} \frac{\partial v}{\partial x_j} dx,$$

and hence

$$(2.1) \quad A(u, v) = \sum_{i=1}^{n_r} A_i(u, v).$$

To define  $B$ , we first decompose functions in  $S_h^0(\Omega)$  as follows: Write  $W = W_P + W_H$  where  $W_P \in S_h^0(\Omega_1) \oplus \cdots \oplus S_h^0(\Omega_{n_r})$  and satisfies, for  $k = 1, \dots, n_r$ ,

$$(2.2) \quad A_k(W_P, \Phi) = A_k(W, \Phi) \quad \text{for all } \Phi \in S_h^0(\Omega_k).$$

Notice that  $W_P$  is determined on  $\Omega_k$  by the values of  $W$  on  $\Omega_k$  and that

$$(2.3) \quad A_k(W_H, \Phi) = 0 \quad \text{for all } \Phi \in S_h^0(\Omega_k).$$

Thus on each  $\Omega_k$ ,  $W$  is decomposed into a function  $W_P$  which vanishes on  $\partial\Omega_k$  and a function  $W_H \in S_h(\Omega_k)$  which satisfies the above homogeneous equations and has the same values as  $W$  on  $\partial\Omega_k$ . We shall refer to such a function  $W_H$  as “discrete  $A$ -harmonic”.

We note that the above decomposition is orthogonal in the inner product defined by  $A$  and hence

$$A(W, W) = A(W_P, W_P) + A(W_H, W_H).$$

We shall define  $B(\cdot, \cdot)$  by replacing the  $A(W_H, W_H)$  term in the above equation. To do this, we first note that by Lemma 3.2 of Part I [4],

$$(2.4) \quad c|W_H|_{1/2, \partial\Omega_k}^2 \leq A_k(W_H, W_H) \leq C|W_H|_{1/2, \partial\Omega_k}^2$$

for discrete  $A$ -harmonic functions  $W_H$  with zero mean value on  $\Omega_k$ . The norm  $|\cdot|_{1/2, \partial\Omega_k}$  is the weighted norm on  $H^{1/2}(\partial\Omega_k)$  given by (see [10], [11])

$$|w|_{1/2, \partial\Omega_k} \equiv \left( \int_{\partial\Omega_k} \int_{\partial\Omega_k} \frac{(w(x) - w(y))^2}{|x - y|} ds(x) ds(y) + d^{-1} \langle w, w \rangle_{\partial\Omega_k} \right)^{1/2}.$$

Here,  $\langle \cdot, \cdot \rangle_{\partial\Omega_k}$  denotes the  $L^2$  inner product on  $\partial\Omega_k$  (the corresponding norm will be denoted by  $|\cdot|_{0,\partial\Omega_k}$ ). In turn, we shall replace the norm  $|\cdot|_{1/2,\partial\Omega_k}$  by a more computationally convenient norm.

To this purpose, we define the operator  $l$  on  $S_h(\partial\Omega_k)$  by

$$(2.5) \quad \langle lV, \Phi \rangle_{\partial\Omega_k} \equiv \langle V', \Phi' \rangle_{\partial\Omega_k} \quad \text{for all } \Phi \in S_h(\partial\Omega_k),$$

where the primes denote differentiation with respect to arc length along each side of  $\partial\Omega_k$ . Now  $l$  is a linear operator on  $S_h(\partial\Omega_k)$  approximating the boundary operator  $-\frac{\partial^2}{\partial s^2}$ , and it can be shown that there are constants  $c$  and  $C$ , independent of  $d$  and  $h$ , such that

$$(2.6) \quad c|V|_{1/2,\partial\Omega_k}^2 \leq \langle l^{1/2}V, V \rangle_{\partial\Omega_k} + d^{-1}\langle V, V \rangle_{\partial\Omega_k} \leq C|V|_{1/2,\partial\Omega_k}^2$$

for all  $V \in S_h(\partial\Omega_k)$ . The following Poincaré inequality holds for all  $W$  with zero mean value on  $\partial\Omega_k$ ,

$$d^{-1}|W|_{0,\partial\Omega_k}^2 \leq cd\langle lW, W \rangle_{\partial\Omega_k}.$$

It then follows by expansion in terms of eigenvectors of  $l$  that

$$|W - \bar{W}|_{1/2,\partial\Omega_k}^2 \leq c\langle l^{1/2}(W - \bar{W}), W - \bar{W} \rangle_{\partial\Omega_k} = c\langle l^{1/2}W, W \rangle_{\partial\Omega_k},$$

where  $\bar{W}$  is the mean value of  $W$  on  $\partial\Omega_k$ . Consequently, we may replace (2.4) by

$$(2.7) \quad c\langle l^{1/2}W_H, W_H \rangle_{\partial\Omega_k} \leq A_k(W_H, W_H) \leq C\langle l^{1/2}W_H, W_H \rangle_{\partial\Omega_k},$$

which holds for all discrete  $A$ -harmonic functions  $W_H$ . Summing the above inequality gives

$$(2.8) \quad c\langle QW_H, W_H \rangle_\Gamma \leq A(W_H, W_H) \leq C\langle QW_H, W_H \rangle_\Gamma,$$

where

$$\langle QW_H, W_H \rangle_\Gamma \equiv \sum_k \alpha_k \langle l^{1/2}W_H, W_H \rangle_{\partial\Omega_k}.$$

The constants  $\alpha_k$  are scaling factors. One reasonable choice is to take  $\alpha_k = (\lambda_1^k + \lambda_0^k)/2$  where  $\lambda_1^k$  and  $\lambda_0^k$  are respectively the largest and smallest eigenvalue of the  $2 \times 2$  matrix  $\{a_{ij}(x_0)\}$  at some point  $x_0 \in \Omega_k$ . By (2.8), the form  $\langle QW_H, W_H \rangle_\Gamma$  is uniformly equivalent to  $A(W_H, W_H)$ . Consequently, the form  $\tilde{B}$  defined by

$$(2.9) \quad \tilde{B}(W, W) \equiv A(W_P, W_P) + \langle QW_H, W_H \rangle_\Gamma$$

is uniformly equivalent to  $A$  on  $S_h^0(\Omega) \times S_h^0(\Omega)$ . The difficulty with using  $\tilde{B}$  as our preconditioner is that the corresponding algorithm for solving (1.5) requires the solution of problems of the form: Find  $V \in S_h(\Gamma)$  such that

$$(2.10) \quad \langle QV, \phi \rangle_\Gamma = F(\phi) \quad \text{for all } \phi \in S_h(\Gamma).$$

It is not easy to solve (2.10); consequently, the choice of  $B = \tilde{B}$  will not lead to a good preconditioner.

Finally, we shall define our preconditioner for  $A$  by replacing  $Q$  in (2.9) by an operator  $\bar{Q}$  which is easier to invert. In fact, we define  $\bar{Q}$  from its inverse. Set

$$(2.11) \quad \bar{Q}^{-1} = P_m(\bar{Q}^{-1}Q)\bar{Q}^{-1},$$

where  $\tilde{Q}$  is some other positive definite symmetric operator on  $S_h(\Gamma)$  and  $P_m$  is a polynomial of degree  $m$ . For our present discussion, we can think of  $\tilde{Q}$  as arbitrary. Some interesting examples, from a computational point of view, will result from the preconditioners constructed in Parts I and II. The possible choices for  $\tilde{Q}$  will be considered in more detail in later sections.

The polynomial in (2.11) is defined (complete details are given in Section 3) so that  $\tilde{Q}^{-1}$  is positive definite and  $\tilde{Q}$  is uniformly equivalent to  $Q$  on  $S_h(\Gamma)$ , i.e.,

$$(2.12) \quad c_0 \langle \tilde{Q}V, V \rangle_\Gamma \leq \langle QV, V \rangle_\Gamma \leq c_1 \langle \tilde{Q}V, V \rangle_\Gamma \quad \text{for all } V \in S_h(\Gamma).$$

Hence, we define our preconditioner  $B$  by

$$(2.13) \quad B(W, W) = A(W_P, W_P) + \langle \tilde{Q}W_H, W_H \rangle_\Gamma.$$

An immediate consequence of (2.12) and (2.8) is that  $B$  is uniformly equivalent to  $A$ ; more precisely, we have the following:

**THEOREM.** *Let  $B$  be given by (2.13), where  $\tilde{Q}$ , defined by (2.11), satisfies (2.12) with  $c_0$  and  $c_1$  independent of  $d$  and  $h$ . Then there exist positive constants  $c$  and  $C$  independent of  $d$  and  $h$  such that*

$$cB(W, W) \leq A(W, W) \leq CB(W, W) \quad \text{for all } W \in S_h^0(\Omega).$$

We shall describe a three-step algorithm to compute the solution  $W = W_P + W_H$  of (1.5) (see [4] and [5]). The function  $W_P$  extended by zero outside of  $\Omega_k$  is a function in  $S_h^0(\Omega_k)$  which satisfies

$$(2.14) \quad A_k(W_P, \Phi) = G(\Phi) \quad \text{for all } \Phi \in S_h^0(\Omega_k).$$

Thus, for step one, the function  $W_P$  on  $\Omega_k$  can be obtained by solving the corresponding Dirichlet problem (2.14). Note that the problems on different subdomains are independent of each other and hence can be solved in parallel.

Now with  $W_P$  known, we are left with the problem of finding  $W_H$ , the second step in the algorithm. It is not difficult to see that the boundary values of  $W_H$  satisfy the equation

$$(2.15) \quad \langle \tilde{Q}W_H, \theta \rangle_\Gamma = G(\bar{\theta}) - A(W_P, \bar{\theta}) \quad \text{for all } \theta \in S_h(\Gamma),$$

where  $\bar{\theta}$  is any extension of  $\theta$  in  $S_h^0(\Omega)$ . Let us discuss the solution of (2.15) in more detail.

To solve (2.15), we must apply the polynomial in (2.11). Let  $V_H$  satisfy  $\tilde{Q}V_H = \tilde{Q}W_H$ , i.e.,

$$(2.16) \quad \langle \tilde{Q}V_H, \theta \rangle_\Gamma = G(\bar{\theta}) - A(W_P, \bar{\theta}) \quad \text{for all } \theta \in S_h(\Gamma).$$

By the definition of  $\tilde{Q}$ ,  $W_H$  on  $\Gamma$  is given by

$$(2.17) \quad W_H = P_m(\tilde{Q}^{-1}\tilde{Q})V_H.$$

In addition, we must evaluate  $\tilde{Q}^{-1}Q$ , i.e., given  $\zeta \in S_h(\Gamma)$  we must find  $\eta = \tilde{Q}^{-1}Q\zeta$  solving

$$(2.18) \quad \langle \tilde{Q}\eta, \theta \rangle_\Gamma = \langle Q\zeta, \theta \rangle_\Gamma \quad \text{for all } \theta \in S_h(\Gamma).$$

Accordingly, the computation of  $W_H$  on  $\Gamma$  only requires evaluation of the form  $\langle Q\zeta, \cdot \rangle_\Gamma$  and the inversion of the  $\langle \tilde{Q}\cdot, \cdot \rangle_\Gamma$  form. The evaluation of the right-hand side of (2.18) is discussed in Section 4.

Once the boundary values of  $W_H$  are known, the third step of the algorithm only requires the computation of the discrete harmonic extension to the interior of the subdomains. As described in [4], [5], this problem can be reduced to the solution of independent Dirichlet problems on the subdomains.

*Remark 2.1.* As will be seen in the following section, the degree of the polynomial  $P_m$  depends upon the relative condition number of the forms  $Q$  and  $\tilde{Q}$ . Indeed, if  $Q$  and  $\tilde{Q}$  satisfy inequalities of the form

$$(2.19) \quad \lambda_0 \langle \tilde{Q}v, v \rangle_\Gamma \leq \langle Qv, v \rangle_\Gamma \leq \lambda_1 \langle \tilde{Q}v, v \rangle_\Gamma \quad \text{for all } v \in S_h(\Gamma),$$

then it suffices to choose  $m$  proportional to  $\sqrt{\lambda_1/\lambda_0}$ .

*Remark 2.2.* Other examples of  $\tilde{Q}$  have been constructed in our earlier papers. If  $\tilde{Q}$  is chosen to be the identity, then (2.19) holds with  $\lambda_1/\lambda_0 \leq c(dh)^{-1}$ . Choosing  $\tilde{Q}$  corresponding to the boundary form constructed in [5], i.e.,

$$(2.20) \quad \langle \tilde{Q}W_H, W_H \rangle_\Gamma \equiv Q(W_H, W_H),$$

where  $Q$  is defined by (2.14) of [5], the results of [5] imply that (2.19) holds with  $\lambda_1/\lambda_0 \leq cd/h$ . Finally, choosing  $\tilde{Q}$  corresponding to the boundary form constructed in [4], i.e.,

$$(2.21) \quad \langle \tilde{Q}W_H, W_H \rangle_\Gamma \equiv B(W_H, W_H),$$

where  $B$  is defined by (2.3) of [4], the results of [4] show that (2.19) holds with  $\lambda_1/\lambda_0 \leq c(1 + \ln(d/h)^2)$ .

**3. The Construction of the Polynomial  $P_m$ .** In this section we shall construct and analyze the polynomial  $P_m$  which appears in (2.11). The ideas involved here are not new, but we will restate the relevant results and constructions for completeness.

We first observe that (2.12) is equivalent to

$$(3.1) \quad c \langle Q^{-1}V, V \rangle_\Gamma \leq \langle \tilde{Q}^{-1}V, V \rangle_\Gamma \leq C \langle Q^{-1}V, V \rangle_\Gamma \quad \text{for all } V \in S_h(\Gamma).$$

Now the operator  $\tilde{Q}^{-1}Q$  is selfadjoint in the inner product given by

$$[u, v] \equiv \langle Qu, v \rangle_\Gamma,$$

and the change of variable  $X = Q^{-1}V$  gives that (3.1) is equivalent to

$$(3.2) \quad c[X, X] \leq [P_m(\tilde{Q}^{-1}Q)\tilde{Q}^{-1}QX, X] \leq C[X, X] \quad \text{for all } X \in S_h(\Gamma).$$

A straightforward spectral argument gives that (3.2) holds (with  $C = 1 + \varepsilon$  and  $c = 1 - \varepsilon$ ) whenever the polynomial  $P_m$  satisfies

$$(3.3) \quad |1 - xP_m(x)| \leq \varepsilon \quad \text{for all } x \in [\lambda_0, \lambda_1],$$

where  $\varepsilon$  is any positive constant less than one not depending on  $d$  or  $h$ , and  $\lambda_0$  and  $\lambda_1$  are the constants appearing in (2.19).

We shall define  $P_m$  in terms of the Chebyshev polynomials. The Chebyshev polynomial  $T_j(y)$  of degree  $j$  is given by

$$T_j(y) = \cos(j \arccos(y)).$$

Define  $P_m$  by

$$(3.4) \quad 1 - xP_m(x) = \frac{T_{m+1}(y(x))}{T_{m+1}(y(0))},$$

where  $y$  is the linear function which takes the interval  $(\lambda_0, \lambda_1)$  into  $(-1, 1)$ , i.e.,

$$y(x) = \frac{2}{\lambda_1 - \lambda_0}x - \frac{\lambda_1 + \lambda_0}{\lambda_1 - \lambda_0}.$$

Since  $|T_{m+1}(y)| \leq 1$  for  $y \in [-1, 1]$ , we have that

$$(3.5) \quad |1 - xP_m(x)| \leq \frac{1}{|T_{m+1}(-\frac{\lambda_1 + \lambda_0}{\lambda_1 - \lambda_0})|} \leq 2 \left( \frac{\sqrt{\gamma} - 1}{\sqrt{\gamma} + 1} \right)^{m+1} \quad \text{for all } x \in [\lambda_0, \lambda_1],$$

where  $\gamma \equiv \lambda_1/\lambda_0$ . The second inequality in (3.5) follows from the identity

$$T_j(y) = \frac{1}{2}[(y + \sqrt{y^2 - 1})^j + (y - \sqrt{y^2 - 1})^j]$$

and elementary manipulations.

To satisfy (3.3), we must choose  $m$  large enough so that

$$(3.6) \quad \left( \frac{\sqrt{\gamma} - 1}{\sqrt{\gamma} + 1} \right)^{m+1} \leq \frac{\varepsilon}{2}.$$

Consequently, it suffices to choose  $m$  in proportion to  $\sqrt{\gamma}$  as  $\gamma$  becomes large. Note however, that for any  $m$ , there is an  $\varepsilon(m, \lambda_1, \lambda_0) < 1$  satisfying (3.3). This implies that  $\tilde{Q}$  will always be positive. The following proposition follows immediately.

**PROPOSITION 3.1.** *Let  $0 < \varepsilon < 1$  be given. There exists a positive constant  $C_1$  independent of  $d$  and  $h$  such that if  $P_m$  is given by (3.4) and*

- (i) *if  $\tilde{Q}$  corresponds to the identity operator on  $S_h(\Gamma)$  and  $m \geq C_1(1 + (dh)^{-1/2})$ ,*  
*or*
- (ii) *if  $\tilde{Q}$  is given by (2.20) and  $m \geq C_1(1 + (d/h)^{1/2})$ , or*
- (iii) *if  $\tilde{Q}$  is given by (2.21) and  $m \geq C_1(1 + \ln(d/h))$ ,*

*then (3.3) holds. Furthermore, (2.12) holds with  $c_0 = 1 - \varepsilon$  and  $c_1 = 1 + \varepsilon$ .*

**Remark 3.1.** To use the preconditioner defined by  $B$  with a given  $\tilde{Q}$ , one needs to know bounds  $\lambda_0$  and  $\lambda_1$  of (2.19). Excellent bounds can be obtained in practice by, for example, applying the power method for eigenvalue estimation. This involves the repeated evaluation of  $\tilde{Q}^{-1}Q$ , which is an essential ingredient in the polynomial evaluation (2.11). The cost of this calculation is minor compared to the overall cost of the algorithm, and no additional coding is necessary.

**Remark 3.2.** The most straightforward algorithm involves choosing  $\tilde{Q}$  to be the identity. For smooth problems,  $\lambda_1/\lambda_0$  is not too large (see Example 5 of Section 5), in which case the degree of  $P_m$  grows like  $(dh)^{-1/2}$ . However, this algorithm requires an excessive number of terms in examples with large jumps in the coefficients across subregion interfaces. In contrast, if we use (2.21) or (2.20) to define  $\tilde{Q}$ , then the constant  $C_1$  appearing in Proposition 3.1 is independent of the jumps in coefficients as long as the jumps occur at the subdomain boundaries (see Examples 2 and 4 of Section 5).

**Remark 3.3.** The coefficients of the polynomial  $P_m$  can easily be calculated by using well-known identities involving Chebyshev polynomials. However, when  $\lambda_1/\lambda_0$  is large, the computation of  $P_m(\tilde{Q}^{-1}Q)$  directly, using the coefficients of  $P_m$ , is somewhat unstable. We suggest the use of the following two-term recurrence relation for  $R_m \equiv P_m(\tilde{Q}^{-1}Q)V$ :

- (i) Define  $\rho = \frac{\lambda_1 - \lambda_0}{\lambda_1 + \lambda_0}$  and  $\alpha = \frac{2}{\lambda_1 + \lambda_0}$ ;



- (ii) set  $\omega_0 = 2$  and  $R_0 = \alpha V$ ;
- (iii) set  $w_1 = (1 - \omega_0 \rho^2/4)^{-1}$  and  $R_1 = \frac{4\alpha}{(2-\rho^2)}[V - \frac{\alpha}{2}\tilde{Q}^{-1}QV]$ ;
- (iv) for  $m > 1$ , set  $\omega_m = (1 - \omega_{m-1} \rho^2/4)^{-1}$  and

$$R_m = \omega_m(R_{m-1} - \alpha\tilde{Q}^{-1}QR_{m-1}) + \alpha\omega_m V - (\omega_m - 1)R_{m-2}.$$

**4. Computational Aspects and Generalizations.** In this section we shall consider various computational aspects of the method as well as some extensions and generalizations. We first describe the computation of the right-hand side of (2.18) in the special case where the mesh points on  $\partial\Omega_k$  are uniformly spaced. We next give a way of extending the techniques of Section 2 to variable coefficient problems on certain irregular mesh domains.

Assume first that the nodes on  $\partial\Omega_j$  are equally spaced with respect to arc length. As discussed in Section 2, given a function  $\zeta \in S_h(\Gamma)$ , we must be able to compute the data

$$\langle Q\zeta, \cdot \rangle_\Gamma$$

appearing in (2.18). By the definition of  $Q$ , it obviously suffices to compute the data

$$(4.1) \quad \langle l^{1/2}\zeta, \cdot \rangle_{\partial\Omega_j}$$

for each subdomain  $\Omega_j$ . We consider first the operator  $l$  from which  $l^{1/2}$  is defined. Let  $r$  be the number of nodes on  $\partial\Omega_j$  and  $\{\Phi_p, p = 1, \dots, r\}$  denote the nodal basis for  $S_h(\partial\Omega_j)$ , where the nodes are listed in, for example, clockwise order. Given the nodal values

$$w \equiv \begin{pmatrix} w_1 \\ \vdots \\ w_r \end{pmatrix}$$

of a function  $W \in S_h(\partial\Omega_j)$ , the nodal values

$$v \equiv \begin{pmatrix} v_1 \\ \vdots \\ v_r \end{pmatrix}$$

of the function  $V = lW$  satisfy

$$Mv = Nw,$$

where

$$(4.2) \quad N_{pq} = \langle l\Phi_p, \Phi_q \rangle_{\partial\Omega_k} \quad \text{and} \quad M_{pq} = \langle \Phi_p, \Phi_q \rangle_{\partial\Omega_k}.$$

In this case of equally spaced nodes, the matrices  $N$  and  $M$  are simultaneously diagonalizable. The eigenvectors are

$$(4.3) \quad \Psi_p = \begin{pmatrix} \exp(\frac{2\pi i p}{r}) \\ \exp(\frac{2\pi i 2p}{r}) \\ \vdots \\ \exp(\frac{2\pi i r p}{r}) \end{pmatrix} \quad \text{for } p = 1, \dots, r.$$

Here  $i$  is the square root of minus one. The corresponding eigenvalues are given by

$$\lambda_p^M = \left( 4 + 2 \cos \left( \frac{2\pi p}{r} \right) \right) \frac{h}{6}$$

and

$$\lambda_p^N = \left( 2 - 2 \cos \left( \frac{2\pi p}{r} \right) \right) / h.$$

It obviously follows that the eigenvalues for the matrix

$$(4.4) \quad L_{pq} = \langle l^{1/2} \Phi_p, \Phi_q \rangle_{\partial \Omega_k}$$

are given by

$$(4.5) \quad \lambda_p^L = \sqrt{\frac{(2 - 2 \cos(\frac{2\pi p}{r}))(4 + 2 \cos(\frac{2\pi p}{r}))}{6}}.$$

Thus one can compute (4.1) by first expanding the nodal values in the basis of eigenvectors (4.3), multiplying by the eigenvalues (4.5) and then computing the nodal values of the resulting expansion. Note that the transformation from nodal values to coordinates in the eigenvector basis (and vice versa) can be computed in  $O(r \ln r)$  operations by use of the fast Fourier transform.

*Remark 4.1.* From the above discussion we see that the amount of work required to evaluate  $\langle Q\zeta, \cdot \rangle_\gamma$  is  $O(\ln(d/h)/dh)$ . For reasonable domain subdivision strategies, the work involved in evaluating  $\tilde{Q}^{-1}$  is also  $O(\ln(d/h)/dh)$ . Thus, the amount of work required for evaluating  $\bar{Q}^{-1}$  is  $O(m \ln(d/h)/dh)$ . This quantity is usually bounded by  $Ch^{-2}$  (see Proposition 3.1).

We next consider the extension of the techniques of Section 2 to the case where the nodes on  $\partial \Omega_k$  are not uniformly spaced (with respect to arc length). Assume that there are  $r$  nodes on  $\partial \Omega_k$ . Let  $R_k$  be a rectangular mesh with a boundary which has an equally spaced mesh of  $r$  nodes. There exists a piecewise linear map  $T_k : \partial \Omega_k \mapsto \partial R_k$  which takes the mesh of  $\partial \Omega_k$  onto that of  $\partial R_k$ . We then define

$$(4.6) \quad \langle \tilde{l}^{1/2} V, V \rangle_{\partial \Omega_k} \equiv \langle \tilde{l}^{1/2} \tilde{V}, \tilde{V} \rangle_{\partial R_k} \quad \text{for all } V \in S_h(\partial \Omega_k),$$

where  $\tilde{V} = V \circ T_k$ . The form  $Q$  is defined by

$$(4.7) \quad \langle QV, V \rangle_\Gamma = \sum_k \alpha_k \langle \tilde{l}^{1/2} V, V \rangle_{\partial \Omega_k}.$$

All of the constructions of Section 2 now go through. Indeed, we decompose  $W = W_P + W_H$  and define  $B$  by (2.13), (2.11) using  $Q$  given by (4.7). The algorithm for computing the solution of (1.5) is completely analogous to that described in Section 2. By (4.6), the evaluation of  $\langle \tilde{l}^{1/2} \zeta, \cdot \rangle_{\partial \Omega_k}$  may be implemented exactly as described in the first part of this section, i.e., we use the procedure given immediately after (4.5).

Finally, we note that in order to get a preconditioner for  $A$ , we can apply our techniques to any other comparable form  $\tilde{A}$ . The form  $\tilde{A}$  is chosen for computational convenience. For example,  $\tilde{A}$  can be chosen so that it can be ‘fast solved’ even when  $A$  corresponds to a variable coefficient operator on a nonuniform mesh as described in Section 4 of [4].

**5. Numerical Experiments.** In this section we shall present some results of numerical experiments which illustrate the convergence properties of the preconditioning methods of this paper. We use (2.13) as a preconditioner in conjunction with the conjugate gradient method. To help illustrate the differences in performance

between the preconditioners of [4], [5] and that of this paper, we will consider the same basic set of examples as those given in Parts I and II. We shall report many of the same parameters as given in Parts I and II. We shall for example, compute the condition number  $K$  of the preconditioned system\*\*. In some examples, we shall also report  $n$ , the number of iterations required to reduce the matrix norm  $(Ax \cdot x)^{1/2}$  of the error  $E_n = U - U_n$  by a specific factor. Here  $U$  is a randomly generated solution of the matrix equations normalized so that  $-1 \leq U \leq 1$  and  $U_n$  is the approximation to  $U$  obtained using  $n$  steps of the iterative algorithm. In addition, we shall include spectral bounds  $K_b$  for the boundary operator  $\tilde{Q}^{-1}Q$  and the degree  $m$  of the polynomial  $P_m$ .

The examples were chosen to illustrate the effectiveness of the algorithm on problems with both smooth and discontinuous coefficients on domains with different geometries. In all of these examples, subspaces  $S_h^0(\Omega)$  of piecewise linear functions defined on a quasi-uniform mesh of size  $h$  were used and the algorithm was applied to solve the finite element equations approximating the solution of an elliptic problem of the form (1.1). The procedure discussed in Section 4 of Part I for choosing the coefficients of the preconditioning form and solving the related subproblems was used throughout this section. In all examples, estimates for the largest and smallest eigenvalue of  $\tilde{Q}^{-1}Q$  were computed by the power method. These estimates were used for  $\lambda_0$  and  $\lambda_1$  in (2.19). The degree  $m$  of the polynomial  $P_m$  was usually taken to be the greatest integer less than or equal to  $1 + \sqrt{K_b}$  where  $K_b \equiv \sqrt{\lambda_1/\lambda_0}$ .

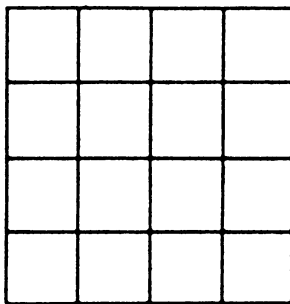


FIGURE 5.1  
*Subdivision of the square.*

*Example 1.* For our first example we take  $L = -\Delta$ , the Laplace operator (i.e.,  $a_{11} = a_{22} = 1$  and  $a_{12} = a_{21} = 0$ ),  $\Omega$  the unit square and a regular rectangular mesh of size  $h$ . Note that, although in this very simple case the resulting equations may be solved rapidly on a serial machine by a variety of ‘fast’ methods, the algorithms of Part I and II would be particularly appealing for a machine with parallel architecture. We will also use this example as a benchmark for the more complicated examples to follow. We subdivide the domain  $\Omega$  into sixteen subregions as indicated in Figure 5.1.

Table 5.1 illustrates the iterative reduction rates for Example 1 when  $h = 1/32$ . The largest and smallest eigenvalues for  $\tilde{Q}^{-1}Q$  were .36 and 1.3, respectively, and  $m$  was taken to be equal to 2. The table lists the total reduction and average reduction

\*\*The condition number  $K$  is defined to be  $\beta_1/\beta_0$  where  $\beta_0$  and  $\beta_1$  are defined to be, respectively, the maximum and minimum constants satisfying (1.6).

rate as a function of the number of iterations in the matrix norm  $(Ax \cdot x)^{1/2}$  and the maximum norm. These reductions are normalized so that the initial error is unity. We see, for example, that a reduction of .0001 in the  $A$ -norm (resp. maximum norm) requires only 7 (resp. 9) iterations.

TABLE 5.1  
*Iterative convergence for Example 1.*

Iteration	$A$ -error	$A$ -error Average Reduction	Max-error	Max-error Average Reduction
1	$1.7 \times 10^{-1}$	.17	$7.8 \times 10^{-1}$	.78
2	$2.9 \times 10^{-2}$	.17	$3.2 \times 10^{-1}$	.57
3	$1.5 \times 10^{-2}$	.24	$1.5 \times 10^{-1}$	.54
4	$3.5 \times 10^{-3}$	.24	$2.7 \times 10^{-2}$	.40
5	$8.0 \times 10^{-4}$	.24	$7.7 \times 10^{-3}$	.38
6	$2.9 \times 10^{-4}$	.26	$2.7 \times 10^{-3}$	.37
7	$8.9 \times 10^{-5}$	.26	$1.3 \times 10^{-3}$	.39
8	$3.9 \times 10^{-5}$	.28	$4.9 \times 10^{-4}$	.39
9	$7.9 \times 10^{-6}$	.27	$6.1 \times 10^{-5}$	.34
10	$1.7 \times 10^{-6}$	.27	$1.6 \times 10^{-5}$	.33
11	$5.6 \times 10^{-7}$	.27	$7.0 \times 10^{-6}$	.34

To more fully illustrate the convergence behavior of the method on this problem, we consider Table 5.2, which gives the condition number and theoretical reduction\*\*\* for Example 1 as a function of the mesh size  $h$ . We note that the theoretical reduction gives a pessimistic bound on the worst case convergence in the  $A$ -norm. For example, the actual reduction rate given in Table 5.1 for 11 iterations was .27, which is considerably better than the theoretical rate of .32 given in Table 5.2 for  $h = 1/32$ .

TABLE 5.2  
*Convergence for Example 1.*

$h$	$K_b$	$m$	$K$	$\rho$	$n$
1/8	1.8	2	2.3	.21	6
1/16	2.6	2	3.0	.27	7
1/32	3.6	2	3.7	.32	7
1/64	5.0	3	3.2	.28	6
1/128	6.2	3	3.5	.30	6

In the next table, we consider the effect that the degree of the polynomial  $P_m$  has on the rate of convergence of the preconditioned algorithm. Table 5.3 gives the number of iterations  $n$  required to reduce the  $A$ -norm error by .0001 and the observed average reduction (in the  $A$ -norm) per iteration as a function of  $m$ . Clearly, as  $m$  tends to infinity, the operator  $\bar{Q}$  tends to  $Q$ . Table 5.3 suggests that the methods converge rapidly, even for small values of  $m$ , and shows that very little

\*\*\*It is well known (cf. [12]) that the error for preconditioned conjugate gradient iteration satisfies  $(AE_n \cdot E_n) \leq 4\rho^{2n}(AE_0 \cdot E_0)$ , where the reduction factor  $\rho$  is given by  $\rho \equiv (\sqrt{K} - 1)/(\sqrt{K} + 1)$ .

$\mu=300$	$\mu=0.0001$	$\mu=31400$	$\mu=5$
$\mu=0.05$	$\mu=8$	$\mu=0.07$	$\mu=2700$
$\mu=10^6$	$\mu=0.1$	$\mu=200$	$\mu=9$
$\mu=1$	$\mu=6000$	$\mu=4$	$\mu=140000$

FIGURE 5.2  
The coefficients for Example 2.

improvement (in the convergence rate of the preconditioned algorithm) results from a more accurate approximation of  $Q^{-1}$ . The results given in the table correspond to  $h = 1/32$ ; similar results were obtained for other values of  $h$ .

TABLE 5.3  
Convergence of the preconditioned algorithm  
as a function of  $m$  for Example 1.

$m$	$K$	Observed Reduction	$n$
1	7.5	.33	9
2	3.7	.27	7
3	2.8	.21	6
4	2.9	.21	6
8	2.8	.21	6

*Example 2.* In this example,  $\Omega$  is the unit square and the subdomains were taken as in Example 1 (see Figure 5.1). The operator  $L$  is taken to have coefficients which have discontinuities across the subdomain boundaries. More specifically, we take  $a_{11} = a_{22} = \mu$  and  $a_{12} = a_{21} = 0$ , where  $\mu$  is the randomly chosen piecewise constant function on the subdomains as indicated in Figure 5.2. Table 5.4 gives the results for the condition number of the preconditioned system and the theoretical reduction factors for this example as a function of  $h$ .

TABLE 5.4  
Convergence results for Example 2.

$h$	$K_b$	$m$	$K$	$\rho$	$n$
1/8	1.9	2	2.3	.21	6
1/16	2.7	2	3.1	.27	6
1/32	3.9	2	3.8	.32	6
1/64	5.0	3	3.4	.29	6
1/128	6.4	3	3.8	.32	6

Note that the results differ only slightly from those given for the Laplacian in Table 5.2. We remark that similar results were obtained in tests with other

randomly chosen coefficients. This indicates that the iterative method of this paper will be extremely effective on interface problems, even when the coefficients change drastically across interfaces, as long as the subdomain boundaries align with the interface boundaries.

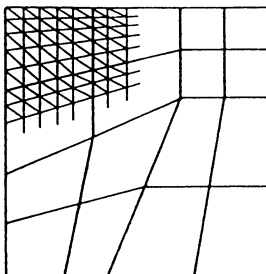


FIGURE 5.3  
*The irregular geometry of Example 3.*

*Example 3.* In this example, we consider an interface problem where the interface separates two domains with irregular geometries. The domain  $\Omega$  is again the unit square subdivided into sixteen subdomains as illustrated in Figure 5.3. The space  $S_h^0(\Omega)$  is taken to be the piecewise linear functions defined on the irregular mesh roughly exemplified by the lighter lines. Again the coefficients of  $L$  are piecewise constant functions defined by  $a_{11} = a_{22} = \mu$  and  $a_{12} = a_{21} = 0$ , where  $\mu$  is given by Figure 5.4.

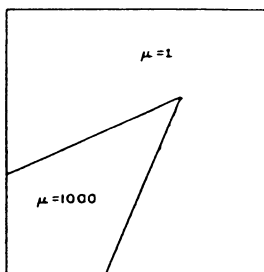


FIGURE 5.4  
*The coefficients of Example 3.*

Results for this problem are given in Table 5.5. A comparison with Table 5.2 indicates that the irregular geometry of this example increased the condition number only by at most a factor of three. This results in less than a factor of two increase in the number of iterations required for a given accuracy. Here again,  $m$  was equal to two or three.

TABLE 5.5  
Convergence results for Example 3.

$h$	$K_b$	$K$	$\rho$	Observed Reduction	$n$
1/8	1.8	4.9	.38	.29	9
1/16	2.6	7.6	.47	.40	11
1/32	3.6	9.9	.52	.45	12
1/64	4.9	8.9	.50	.42	11

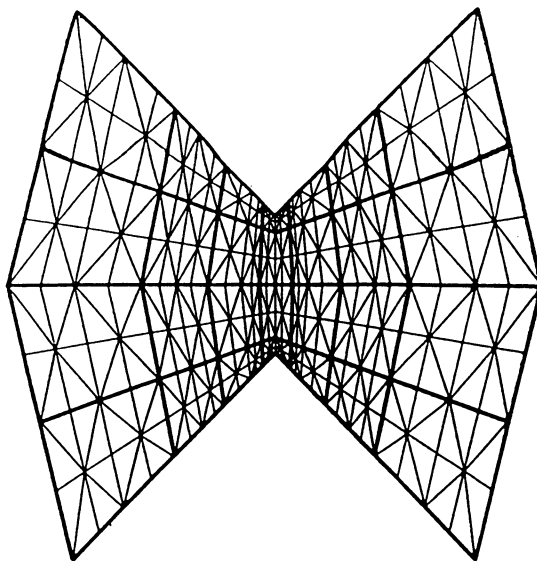


FIGURE 5.5  
The mesh and subdomain structure for Example 4.

*Example 4.* In this example, we illustrate the present algorithm applied to the solution of a problem on a polygonal domain with nonconvex corners. The mesh and subdomain structure were chosen as illustrated in Figure 5.5. Note the mild refinement near the nonconvex corners of the domain. For the operator  $L$  we use the Laplacian as in Example 1. The results for this case are given in Table 5.6.

TABLE 5.6  
Convergence results for Example 4.

Number of Unknowns	$K_b$	$K$	$\rho$	Observed Reduction	$n$
405	2.8	4.4	.35	.35	9
1705	3.8	5.8	.41	.40	10
6993	5.2	5.3	.40	.37	10

*Example 5.* As a final example, we illustrate the algorithm described in Remark 3.2, i.e., we consider the case where  $\tilde{Q} = I$ . We consider the problem and domain decomposition of Example 2. Table 5.7 gives the condition number  $K_b$  of  $\tilde{Q}$ ,  $n$ ,  $K$ , and the observed reduction in the  $A$ -norm as a function of  $h$ . In this case, we increased  $m$  as suggested by Proposition 3.1 (i).

TABLE 5.7  
Convergence results for Example 5.

$h$	$K_b$	$m$	$K$	Observed Reduction	$n$
1/8	11.5	4	2.4	.21	6
1/16	25	5	3.2	.25	7
1/32	52	8	3.3	.31	7
1/64	105	11	4.3	.31	8

Department of Mathematics  
Cornell University  
Ithaca, New York 14853  
E-mail: bramble@mathvax.msi.cornell.edu

Brookhaven National Laboratory  
Upton, New York 11973  
E-mail: pasciak@bnl.arpa

Department of Mathematics  
Cornell University  
Ithaca, New York 14853

1. P. E. BJØRSTAD & O. B. WIDLUND, "Solving elliptic problems on regions partitioned into substructures," *Elliptic Problem Solvers II* (G. Birkhoff and A. Schoenstadt, eds.), Academic Press, New York, 1984, pp. 245–256.

2. P. E. BJØRSTAD & O. B. WIDLUND, "Iterative methods for the solution of elliptic problems on regions partitioned into substructures," *SIAM J. Numer. Anal.*, v. 23, 1986, pp. 1097–1120.

3. J. H. BRAMBLE, J. E. PASCIAK & A. H. SCHATZ, "An iterative method for elliptic problems on regions partitioned into substructures," *Math. Comp.*, v. 46, 1986, pp. 361–369.

4. J. H. BRAMBLE, J. E. PASCIAK & A. H. SCHATZ, "The construction of preconditioners for elliptic problems by substructuring, I," *Math. Comp.*, v. 47, 1986, pp. 103–134.

5. J. H. BRAMBLE, J. E. PASCIAK & A. H. SCHATZ, "The construction of preconditioners for elliptic problems by substructuring, II," *Math. Comp.*, v. 49, 1987, pp. 1–16.

6. B. L. BUZBEE & F. W. DORR, "The direct solution of the biharmonic equation on rectangular regions and the Poisson equation on irregular regions," *SIAM J. Numer. Anal.*, v. 11, 1974, pp. 753–763.

7. B. L. BUZBEE, F. W. DORR, J. A. GEORGE & G. H. GOLUB, "The direct solution of the discrete Poisson equation on irregular regions," *SIAM J. Numer. Anal.*, v. 8, 1971, pp. 722–736.

8. Q. V. DIHN, R. GLOWINSKI & J. PÉRIAUX, "Solving elliptic problems by domain decomposition methods," in *Elliptic Problem Solvers II* (G. Birkhoff and A. Schoenstadt, eds.), Academic Press, New York, 1984, pp. 395–426.

9. G. H. GOLUB & D. MEYERS, *The Use of Preconditioning Over Irregular Regions*, Proc. 6th Internat. Conf. Comput. Meth. Sci. and Engrg., Versailles, 1983.

10. J. L. LIONS & E. MAGENES, *Problèmes aux Limites non Homogènes et Applications*, Dunod, Paris, 1968.

11. J. NEČAS, *Les Méthodes Directes en Théorie des Équations Elliptiques*, Academia, Prague, 1967.

12. W. M. PATTERSON, 3rd, *Iterative Methods for the Solution of a Linear Operator Equation in Hilbert Space—A Survey*, Lecture Notes in Math., vol. 394, Springer-Verlag, New York, 1974.