

Finite Element Methods Are Not Always Optimal

ARTHUR G. WERSCHULZ

Division of Science and Mathematics
Fordham University / College at Lincoln Center
113 West 60th Street
New York, NY 10023

Department of Computer Science
Columbia University
New York, NY 10027

CUCS-178-84

June, 1984

ABSTRACT. Consider a regularly-elliptic $2m$ -th order boundary problem $Lu = f$ with $f \in H^r(\Omega)$, $r \geq -m$. In previous work, we showed that the finite-element method (FEM) using piecewise polynomials of degree k is asymptotically optimal when $k \geq 2m - 1 + r$. In this paper, we show that the FEM is not asymptotically optimal when this inequality is violated. However, there exists an algorithm, the Traub-Wasilkowski-Woźniakowski spline algorithm, which is always optimal. Moreover, the error of the FEM can be arbitrarily larger than the error of the spline algorithm. We also obtain a necessary and sufficient condition for a Galerkin method (or a generalized Galerkin method) to be a spline algorithm.

Keywords and phrases. Elliptic problems, variational methods, finite element methods, optimal algorithms, computational complexity.

1980 Mathematics subject classifications (Amer. Math. Soc.): Primary: 65N30, 68C25. Secondary: 65J10, 65N05, 65N15.

Division of Science and Mathematics, Fordham University / College at Lincoln Center, 113 West 60th Street, New York, NY 10023

Department of Computer Science, Columbia University, New York, NY 10027

This research was supported in part by the National Science Foundation under Grants MCS-8203271 and MCS-8303111.

1. Introduction

This paper deals with the optimal solution of $2m$ -th order regularly-elliptic boundary-value problems $Lu = f$ with $f \in H^r(\Omega)$, $\Omega \subseteq \mathbb{R}^N$. We consider the variational form of such problems having homogeneous boundary conditions (see Section 2). We wish to solve such problems using information of cardinality at most n . (In this Introduction, we have to use words such as information, cardinality, algorithm, etc., without definition; they are defined rigorously in Section 3.)

In [10], we showed that the optimal energy-norm error of an algorithm using information of cardinality n is $\Theta(n^{-(m+r)/N})$, as $n \rightarrow \infty$.^{*} Moreover, this optimal error is achieved by a finite-element method (FEM) using piecewise polynomials of degree k , where $k \geq 2m - 1 + r$. Suppose that this inequality is violated. For instance, suppose we have a program written using piecewise-linear polynomials to solve a second-order problem in a planar region Ω . For $f \in L_2(\Omega)$, this FEM has H^1 -error $\Theta(n^{-1/2})$, which is optimal. What happens when we use this program to solve a problem with (say) $f \in H^1(\Omega)$? Does the error of the FEM improve when f has additional smoothness, and if not, is there a method using the same information as the FEM, but with better error?

^{*} Here and in what follows, we use the Ω - and Θ -notations of Knuth [5], as well as the usual O -notation. That is,

$$f = \Omega(g) \quad \text{if} \quad g = O(f)$$

and

$$f = \Theta(g) \quad \text{if} \quad f = O(g) \quad \text{and} \quad f = \Omega(g).$$

In Section 4, we show that the error of the FEM is $\Theta(n^{-\mu}/N)$ as $n \rightarrow \infty$, where $\mu = \min(k + 1 - m, m + r)$, so that $k \geq 2m - 1 + r$ is a necessary and sufficient condition for the FEM to be asymptotically optimal. On the other hand, we analyze the Traub-Wasilkowski-Wozniakowski spline algorithm (see [9]) which uses the same information as the FEM. We show that the error of the spline algorithm is $\Theta(n^{-(m+r)}/N)$ as $n \rightarrow \infty$, regardless of whether $k \geq 2m - 1 + r$; it is therefore always asymptotically optimal. Moreover, unlike the FEM, the proof of the error estimate for the spline algorithm does not require the "shift theorem"; hence, the spline algorithm is applicable to a wider range of problems.

The optimality result mentioned above is for a worst-case f . Also of interest is the local error, i.e., the error for any particular f . The spline algorithm is known to be strongly optimal, that is, it enjoys optimal local error. It is well-known that the FEM is a Galerkin method. Furthermore, as we show in Section 5, the spline algorithm is a generalized Galerkin method. (Indeed, given the spline algorithm, we show how to realize it as a generalized Galerkin method.) This motivates our interest in the local error of generalized Galerkin methods.

The ratio of the local error of an algorithm to the optimal local error is called the deviation of the algorithm. We pose and solve the question of when a generalized Galerkin method has finite deviation. We show the deviation is finite if and only if the generalized Galerkin method is the spline algorithm.

Do FEM's always have finite deviation? We show the answer is no, by exhibiting an FEM which is not a spline algorithm. We conjecture that no convergent FEM has finite deviation.

In Section 6, we discuss the complexity of obtaining ϵ -approximations. We show that the penalty for using the FEM when $k < 2m - 1 + r$ is unbounded as $\epsilon \rightarrow 0$. Since this is an asymptotic measure, it is useful to know whether the spline algorithm has lower complexity than the FEM for moderate values of ϵ . We show that this is indeed the case, by exhibiting a model problem for which the spline algorithm has lower complexity than the FEM for any ϵ roughly less than one-half.

Finally, in Section 7, we briefly discuss implementation of the spline algorithm, and ask whether it is practical to use.

2. The Variational Boundary-Value Problem

In what follows, we use the standard notations for Sobolev spaces, inner products, and norms, multi-indices, etc. found in Ciarlet [2]. Fractional- and negative-order Sobolev spaces are defined by Hilbert-space interpolation and duality, respectively; see Chapter 2 of [1] and Chapter 4 of [6] for details.

Let Ω be a bounded C^∞ region in \mathbb{R}^N . Define the properly elliptic operator

$$(2.1) \quad Lv := \sum_{|\alpha|, |\beta| \leq m} (-1)^{|\alpha|} D^\alpha (a_{\alpha\beta} D^\beta v)$$

(with real coefficients $a_{\alpha\beta} \in C^\infty(\bar{\Omega})$ such that $a_{\alpha\beta} = a_{\beta\alpha}$) and a normal family of operators

$$(2.2) \quad B_j v := \sum_{|\alpha| \leq q_j} b_{j\alpha} D^\alpha v \quad (0 \leq j \leq m-1),$$

(with real coefficients $b_{j\alpha} \in C^\infty(\partial\Omega)$), where

$$(2.3) \quad 0 \leq q_0 \leq q_1 \leq \dots \leq q_{m-1} \leq 2m-1,$$

which covers L on $\partial\Omega$. Setting

$$(2.4) \quad m^* := \min\{j : q_j \geq m\},$$

we additionally assume that

$$(2.5) \quad \{q_j\}_{j=0}^{m^*-1} \cup \{2m-1 - q_j\}_{j=m^*}^{m-1} = \{0, \dots, m-1\}.$$

(See Chapter 3 of [1], Chapter 5 of [6] for further definitions and illustrative examples.)

Let

$$(2.6) \quad H_E^m(\Omega) := \{v \in H^m(\Omega) : B_j v = 0 \text{ for } 0 \leq j \leq m^* - 1\}$$

denote the space of $H^m(\Omega)$ -functions satisfying the essential boundary conditions of order at most $m^* - 1$. We define a symmetric, continuous bilinear form B on $H_E^m(\Omega)$ by

$$(2.7) \quad B(v, w) := \sum_{|\alpha|, |\beta| \leq m} \int_{\Omega} a_{\alpha\beta} D^{\alpha} v D^{\beta} w.$$

We additionally assume that B is $H_E^m(\Omega)$ -coercive, so that B is an inner-product on $H_E^m(\Omega)$, yielding a norm $\|\cdot\|_B$ defined by

$$(2.8) \quad \|v\|_B := B(v, v)^{1/2}$$

which is equivalent to the norm $\|\cdot\|_m$ on $H_E^m(\Omega)$.

We now define the variational boundary-value problem as follows. Let $r \geq -m$. Given $f \in H^r(\Omega)$, find $u = Sf \in H_E^m(\Omega)$ satisfying

$$(2.9) \quad B(u, v) = (f, v)_0 = \int_{\Omega} f v \quad \forall v \in H_E^m(\Omega).$$

From the Lax-Milgram lemma, S is a Hilbert space isomorphism of $H^{-m}(\Omega)$ onto $H_E^m(\Omega)$, so that $S : H^r(\Omega) \rightarrow H_E^m(\Omega)$ is a bounded linear operator.

It is useful to recall the "shift theorem" (Chapter 3 of [1], Chapter 5 of [6]) which states that since $f \in H^r(\Omega)$, we have $Sf \in H_E^m(\Omega) \cap H^{2m+r}(\Omega)$; moreover, there exists a positive constant σ , independent of f , such that

$$(2.10) \quad \sigma^{-1} \|Sf\|_{2m+r} \leq \|f\|_r \leq \sigma \|Sf\|_{2m+r}.$$

If $r > N/2$, the shift theorem, Sobolev's embedding theorem, and an m -fold integration by parts yield that $u = Sf$ is a classical solution to the problem of finding $u : \bar{\Omega} \rightarrow \mathbb{R}$ satisfying

$$(2.11) \quad \begin{aligned} Lu &= f && \text{in } \Omega \\ B_j u &= 0 && \text{on } \partial\Omega \quad (0 \leq j \leq m-1). \end{aligned}$$

3. Information and Algorithms

In this section, we define a number of the concepts mentioned in the Introduction. Most of the terminology and results are from [9]. As we state these definitions and results, we will illustrate them for the finite element method and Galerkin information.

Recall that we are trying to approximate the transformation $S : H^r(\Omega) \rightarrow H_E^m(\Omega)$ with $r \geq -m$. Since S is not of finite rank and we wish to use finite algorithms, we are only allowed to sample a finite amount of information about problem elements $f \in H^r(\Omega)$. Here (linear) information of cardinality n is a surjective linear mapping $n : H^r(\Omega) \rightarrow \mathbb{R}^n$, so that we may write

$$(3.1) \quad nf = [\lambda_1(f) \ \dots \ \lambda_n(f)]^T \quad \forall f \in H^r(\Omega)$$

where $\lambda_1, \dots, \lambda_n$ are linearly independent linear functionals on $H^r(\Omega)$. (See Chapter 7 of [9] for a discussion of why we consider only linear information.)

Example 3.1. Let \mathcal{S} be a subspace of $H_E^m(\Omega)$ of dimension n and having basis $\{s_1, \dots, s_n\}$. Define $n_{\mathcal{S}} : H^r(\Omega) \rightarrow \mathbb{R}^n$ by

$$(3.2) \quad n_{\mathcal{S}} f := [(f, s_1)_0 \ \dots \ (f, s_n)_0]^T \quad \forall f \in H^r(\Omega).$$

Then $n_{\mathcal{S}}$ is linear information of cardinality n . (Conversely, given any linear information n of cardinality n , one can show that there exists a subspace $\mathcal{S} \subset H_E^m(\Omega)$ of dimension n such that $n = n_{\mathcal{S}}$ if and only if n has an extension to all of $H^{-m}(\Omega)$ which is bounded in the $\|\cdot\|_{-m}$ norm.) We call $n_{\mathcal{S}}$ the Galerkin information generated by \mathcal{S} . ■

In the remainder of this paper, for any Hilbert space H , we denote the unit ball of H by BH , i.e.,

$$(3.2) \quad BH := \{f \in H : \|f\|_H \leq 1\}.$$

In particular, we will let

$$(3.3) \quad \mathcal{F}_0 := BH^F(\Omega).$$

By an algorithm φ using n , we mean a (possibly nonlinear) mapping $\varphi : D_\varphi \subset n(\mathcal{F}_0) \rightarrow H_E^m(\Omega)$. The (worst-case) error $e(\varphi)$ of φ is given by

$$(3.4) \quad e(\varphi) := \sup_{f \in \mathcal{F}_0} \|Sf - \varphi(nf)\|_B.$$

(The restriction to $f \in \mathcal{F}_0$, rather than considering the sup over all $f \in H^F(\Omega)$, is a normalization which is necessary for the error to be finite.) We use the norm $\|\cdot\|_B$ rather than the equivalent norm $\|\cdot\|_m$ for technical reasons, as illustrated in

Example 3.1 (continued). Define the Galerkin method φ_g using n_g by

$$(3.5) \quad \varphi_g(n_g f) := u_g,$$

where $u_g \in \mathcal{S}$ satisfies

$$(3.6) \quad B(u_g, s) = (f, s)_0 \quad \forall s \in \mathcal{S}.$$

Then standard results ([1],[2],[7]) yield

$$(3.7) \quad e(\varphi_g) = \sup_{f \in \mathcal{F}_0} \inf_{s \in \mathcal{S}} \|Sf - s\|_B.$$

In particular, let $\mathcal{S} = \mathcal{S}_n$, where $\{\mathcal{S}_n\}_{n=1}^{\infty}$ is a regular family of finite element subspaces of degree k , i.e., \mathcal{S}_n is an n -dimensional subspace of $H_E^m(\Omega)$ consisting of piecewise polynomials of degree k over a triangulation \mathcal{T}_n of Ω . Here, $\{\mathcal{T}_n\}_{n=1}^{\infty}$ is regular in the sense of page 132 of [2], which (roughly) means that the subregions do not become geometrically degenerate, and that their diameters tend to zero as $n \rightarrow \infty$. (Of course, since Ω is C^∞ , we must make an additional assumption about boundary elements to guarantee that $\mathcal{S}_n \subset H_E^m(\Omega)$; for instance, we may decide to use curved elements as in [3].)

For this case, we write φ_n^* and n_n^* rather than $\varphi_{\mathcal{S}}$ and $n_{\mathcal{S}}$, calling φ_n^* the finite element method (FEM) using \mathcal{S}_n . Suppose now that $\{\mathcal{T}_n\}_{n=1}^{\infty}$ is quasi-uniform (see pg. 272 of [6]), which means that the ratio of the diameters of any two subregions in \mathcal{T}_n is bounded, independent of n . Then the standard results ([1],[2],[6]) yield

$$(3.8) \quad e(\varphi_n^*) = O(n^{-\mu/N}) \text{ as } n \rightarrow \infty \quad \mu = \min(k+1-m, m+r).$$

Moreover, results of Strang and Fix [8] imply that the "O" may be changed to " Θ " when $H_E^m(\Omega) = H^m(\Omega)$ (i.e., no essential boundary conditions), the triangulations \mathcal{T}_n are uniform, and $k \leq 2m - 1 + r$. In Section 4 of this paper, we will remove these three restrictions, so that the bound (3.8) is always sharp. ■

Given information n of cardinality n , we wish to find the minimum error of an algorithm φ using n . In order to do this, let

$$(3.9) \quad \forall f := \{\tilde{f} \in \mathcal{F}_0 : n\tilde{f} = nf\} \quad \forall f \in \mathcal{F}_0.$$

If φ uses n , then knowing only nf , it is impossible for φ to determine which of the elements of the set

$$(3.10) \quad Uf := SVf$$

is being approximated, so that

$$(3.11) \quad e(\varphi) = \sup_{f \in \mathcal{F}_0} e(\varphi, f)$$

where the local error $e(\varphi, f)$ is given by

$$(3.12) \quad e(\varphi, f) := \sup_{\tilde{f} \in V(f)} \|S\tilde{f} - \varphi(nf)\|_B \quad \forall f \in \mathcal{F}_0.$$

Define the local radius $\text{rad } Uf$ by

$$(3.13) \quad \text{rad } Uf := \inf_{a \in H_E^m(\Omega)} \sup_{\tilde{f} \in V(f)} \|a - S\tilde{f}\|_B \quad \forall f \in \mathcal{F}_0.$$

As in Chapter 1 of [9], we have

$$(3.14) \quad \inf_{\varphi} e(\varphi, f) = \text{rad } Uf,$$

so that

$$(3.15) \quad \inf_{\varphi} e(\varphi) = r(n) := \sup_{f \in \mathcal{F}_0} \text{rad } Uf,$$

where $r(n)$ is called the radius of information n . In our Hilbert space setting, one can show that

$$(3.16) \quad r(n) = \sup_{z \in \mathcal{F}_0 \cap \ker n} \|Sz\|_B$$

(see Chapter 2 of [9]).

Example 3.1 (continued). In Section 4, we will show that for quasi-uniform $\{\tau_n\}_{n=1}^{\infty}$,

$$(3.17) \quad r(n_n^*) = \Theta(n^{-(m+r)/N}) \text{ as } n \rightarrow \infty.$$

Hence the FEM φ_n^* has (asymptotically) optimal error using n_n (as $n \rightarrow \infty$) if and only if $k \geq 2m - 1 + r$. ■

Remark 3.1. Now that we know the minimal local and worst-case errors of algorithms using n , it is useful to find algorithms achieving these minima. Let $P : H^r(\Omega) \rightarrow H^r(\Omega)$ denote the orthogonal projector onto $(\ker n)^\perp$. Define the (Traub-Wasilkowski-Woźniakowski) spline algorithm φ^S (Chapter 4 of [9]) by

$$(3.18) \quad \varphi^S(nf) := SPf \quad \forall f \in \mathcal{F}_0.$$

One may check that φ^S is well-defined, and that

$$(3.19) \quad e(\varphi^S, f) = \text{rad } Uf \quad \forall f \in \mathcal{F}_0,$$

which implies that φ^S is an optimal error algorithm, i.e., for any φ using n ,

$$(3.20) \quad e(\varphi^S) = r(n) \leq e(\varphi).$$

Not only is φ^S an optimal error algorithm, but it is a strongly optimal error algorithm, that is,

$$(3.21) \quad e(\varphi^S, f) \leq e(\varphi, f) \quad \forall f \in \mathcal{F}_0.$$

We will discuss FEM's and spline algorithms more fully in Section 5. ■

Just as we can ask which algorithm makes optimal use of given information, one can ask which information of a given cardinality is best. Let

$$(3.22) \quad r(n) := \inf\{r(\alpha) : \alpha \text{ is of cardinality at most } n\}$$

denote the n th minimal radius of information; we say that α of cardinality at most n is an n th optimal information if

$$(3.23) \quad r(\alpha) = r(n).$$

Then (Chapter 2 of [9])

$$(3.24) \quad r(n) = d_n(S(\mathcal{F}_0), H_E^m(\Omega)),$$

where the Kolmogorov n -width of a balanced subset X of a Hilbert space H with norm $\|\cdot\|_H$ is given by

$$(3.25) \quad d_n(X, H) := \inf\{\sup_{x \in X} \inf_{y \in A_n} \|x - y\|_H : A_n \text{ subspace of } H, \dim A_n \leq n\}.$$

Example 3.1 (continued). Results from [10] yield that

$$(3.26) \quad r(n) = \Theta(n^{-(r+m)/N}) \text{ as } n \rightarrow \infty.$$

Hence, the results in Section 4 will imply that α_n^* is (asymptotically, as $n \rightarrow \infty$) an n th optimal information. ■

4. Worst-Case Error of the FEM and the Spline Algorithm

In this section, we show that $k \geq 2m - 1 + r$ is a necessary and sufficient condition for an FEM to have optimal error to within a constant, independent of n . We also show that the spline method is an optimal error algorithm using the n th optimal (to within a constant) information n_n^* , regardless of whether $k \geq 2m - 1 + r$.

Recall that S_n is an n -dimensional subspace of $H_E^m(\Omega)$ consisting of piecewise polynomials of degree k from a triangulation T_n of Ω . We first show

Lemma 4.1. $k \geq m$.

Proof: Suppose on the contrary that $k \leq m - 1$. Since, for any $s \in S_n$, $s \in H^m(\Omega)$ and $s|_K \in C^0(K)$ for each $K \in T_n$, an obvious extension of Theorem 4.2.1 of Ciarlet [2] yields that $S_n \subset C^{m-1}(\Omega)$. Choose $s \in S_n$. Let K_1, K_2 be adjacent elements in the triangulation, let

$$(4.1) \quad F := \partial K_1 \cap \partial K_2$$

and let

$$(4.2) \quad s_i := s|_{K_i} \quad (i = 1, 2).$$

Let $s^* \in P_k$ satisfy $s^* = s_1$ on K_1 ; that is, s^* is s_1 , but treated as a polynomial over Ω rather than over K_1 . Pick a point p on F , and draw a ξ -axis G_p perpendicular to F through p . Hence there is an affine transformation $F_p : \mathbb{R} \rightarrow \mathbb{R}^N$ which is a bijection of \mathbb{R} onto G_p , such that $p = F_p(0)$.

For $i = 1$ and $i = 2$, define

$$(4.3) \quad \sigma_i(\xi) := s_i(F_p(\xi))$$

so that

$$(4.4) \quad \sigma_i(\xi) = s_i(x) \text{ for } x = F_p(\xi).$$

Since

$$(4.5) \quad s_i \in P_k \text{ and } F_p \text{ is affine.}$$

we see that

$$(4.6) \quad \sigma_i(\xi) \text{ is a polynomial of degree at most } k \text{ in } \xi.$$

On the other hand, since $s \in C^{m-1}(\Omega)$, we must have

$$(4.7) \quad \sigma_1^{(j)}(0) = \sigma_2^{(j)}(0) \quad (0 \leq j \leq m-1).$$

Using (4.6), (4.7), and $k \leq m-1$, we see that

$$(4.8) \quad \sigma_1(\xi) = \sigma_2(\xi) \quad \forall \xi \in \mathbb{R},$$

i.e.,

$$(4.9) \quad s_1(x) = s_2(x) \quad \forall x \in G_p \cap (K_1 \cup K_2).$$

Since $p \in F$ is arbitrary, we let p vary along F to find

$$(4.10) \quad s_1(x) = s_2(x) \quad \forall x \in K_1 \cup K_2,$$

i.e.,

$$(4.11) \quad s|_{K_i}(x) = s^*(x) \quad \forall x \in K_1 \cup K_2, \quad i = 1, 2.$$

Repeating this argument, we see that for any $K \in \mathcal{S}_n$,

$$(4.12) \quad s|_K = s^*|_K,$$

and so $s = s^* \in P_K$. Thus s arbitrary in \mathcal{S}_n yields

$$(4.13) \quad \mathcal{S}_n \subseteq P_K,$$

so that

$$(4.14) \quad n = \dim \mathcal{S}_n \leq \dim P_K = \binom{k+N}{N},$$

which is impossible, since k and N are fixed, while n is an arbitrary positive integer. Hence $k \geq m$. ■

We are now able to establish the sharpness of the usual estimate for this error of the FEM φ_n^* , generalizing the work of Strang and Fix [8].

Theorem 4.1. Let $r \geq -m$, and define

$$(4.15) \quad \mu = \min(k+1-m, m+r).$$

Then

$$(i) \quad e(\varphi_n^*) = \Omega(n^{-\mu}/N) \text{ as } n \rightarrow \infty,$$

and

$$(ii) \quad e(\varphi_n^*) = \Theta(n^{-\mu}/N) \text{ as } n \rightarrow \infty \text{ for quasi-uniform } \{\mathcal{S}_n\}_{n=1}^{\infty}.$$

Proof: First note that (3.15), (3.22), (3.23), and (3.26) yield

$$(4.16) \quad e(\varphi_n^*) \geq r(n) = \Theta(n^{-(r+m)}/N).$$

It remains to show

$$(4.17) \quad e(\varphi_n^*) = \Omega(n^{-(k+1-m)/N}) \text{ as } n \rightarrow \infty,$$

since (4.15), (4.16), and (4.17) imply (i), while (i) and the usual estimate (3.8) yield (ii).

In order to show (4.17), we will rely heavily on the notation found in [2]. First, let Ω° be the interior of a hypercube such that $\overline{\Omega^{\circ}} \subset \Omega$,

$$(4.18) \quad \mathcal{J}_n^{\circ} := \{K \in \mathcal{J}_n : K \subset \overline{\Omega^{\circ}}\},$$

and

$$(4.19) \quad \Omega_n := \text{int } \cup \{K : K \in \mathcal{J}_n^{\circ}\}.$$

For any element $K \in \mathcal{J}_n$, let

$$(4.20) \quad \rho_K := \sup\{\text{diam}(S) : S \text{ a ball}, S \subset K\}$$

and

$$(4.21) \quad h_K := \text{diam } K.$$

Then $\{\mathcal{J}_n\}_{n=1}^{\infty}$ regular means that

$$(4.22) \quad \lim_{n \rightarrow \infty} \sup_{K \in \mathcal{J}_n} h_K = 0$$

and there is a constant $\sigma > 0$ such that

$$(4.23) \quad h_K \leq \sigma \rho_K \quad \forall K \in \mathcal{J}_n, \quad \forall n \geq 1.$$

Using (4.18), (4.19), (4.21), and (4.22), we find that

$$(4.24) \quad \bar{\Omega}_n \subset \bar{\Omega}^0 \quad \forall n \geq 1,$$

but that

$$(4.25) \quad \lim_{n \rightarrow \infty} \text{vol}(\bar{\Omega}_n) = \text{vol}(\bar{\Omega}^0).$$

We next choose u to be any function in $H_E^m(\Omega) \cap H^{2m+r}(\Omega)$ such that

$$(4.26) \quad u(x) = \frac{1}{(k+1)!} x_1^{k+1} \quad \forall x \in \bar{\Omega}^0.$$

Let $K \in \mathcal{J}_n^0$. We now claim that there is a constant $C_1 > 0$, independent of K and n , such that

$$(4.27) \quad \inf_{s \in P_k(K)} |u - s|_{m,K}^2 \geq C_1^2 \text{vol}(K)^{2(k+1-m)/N+1},$$

$P_k(K)$ denoting polynomials of degree k over K . To show (4.27), there is an affine bijection $F_K : \hat{K} \rightarrow K$ with

$$(4.28) \quad F_K \hat{x} = B_K \hat{x} + b_K,$$

where \hat{K} is a reference element independent of K , so that K is the F_K -image of a "reference element" \hat{K} independent of n and K . Then Theorem 3.1.2 of [2] yields the existence of a constant $c_1 = c_1(k,m) > 0$ such that

$$(4.29) \quad \inf_{s \in P_k(K)} |u - s|_{m,K}^2 \geq c_1 |\det B_K| \|B_K\|^{-2m} \inf_{\hat{s} \in P_k(\hat{K})} |\hat{u} - \hat{s}|_{m,\hat{K}}^2,$$

where $\|\cdot\|$ is the Euclidean matrix norm and where, for any function $v : K \rightarrow \mathbb{R}$, we define $\hat{v} : \hat{K} \rightarrow \mathbb{R}$ by

$$(4.30) \quad \hat{v}(\hat{x}) := v(x) \text{ for } x = F_K \hat{x}.$$

On the other hand, one may check that the functionals

$$(4.31) \quad \hat{v} \mapsto |\hat{v}|_{k+1, \hat{K}}$$

and

$$(4.32) \quad \hat{v} \mapsto \inf_{\hat{s} \in P_k(\hat{K})} |\hat{v} - \hat{s}|_{m, \hat{K}}$$

are seminorms on $P_{k+1}(\hat{K})$. Since $k \geq m$, we find that these two seminorms have the same kernel, namely the space $P_k(\hat{K})$. Since $P_{k+1}(\hat{K})$ is finite dimensional, there is a constant $c_2 = c_2(k, m, \hat{\Omega}) > 0$ such that

$$(4.33) \quad \inf_{\hat{s} \in P_k(\hat{K})} |\hat{v} - \hat{s}|_{m, \hat{K}} \geq c_2 |\hat{v}|_{k+1, \hat{K}} \quad \forall \hat{v} \in P_{k+1}(\hat{K}).$$

Hence, we may use (4.27) and (4.33) with \hat{v} replaced by \hat{u} to see that

$$(4.34) \quad \inf_{s \in P_k(K)} |u - s|_{m, K}^2 \geq c_1 c_2^2 |\det B_K| \|B_K\|^{-2m} |\hat{u}|_{k+1, \hat{K}}^2$$

$$\geq c_3 \frac{(\|B_K^{-1}\|^{-1})^{2(k+1)}}{\|B_K\|^{2m}} |u|_{k+1, K}^2$$

for $c_3 = c_3(k, m) > 0$ by Theorem 3.1.2 of [2]. Using (4.23) and Theorem 3.1.3 of [2], we find that there exist $c_4, c_5 > 0$ such that

$$(4.35) \quad \begin{aligned} \|B_K^{-1}\|^{-1} &\geq c_4 \rho_K \geq c_4 \sigma^{-1} h_K \\ \|B_K\| &\leq c_5 h_K \end{aligned}$$

Since

$$(4.36) \quad \text{vol}(K) \leq \sigma_N h_K^N,$$

where σ_N is the volume of the unit ball in \mathbb{R}^N , we have

$$(4.37) \quad \inf_{s \in P_k(K)} |u - s|_{m,K}^2 \geq C_1^2 \text{vol}(K)^{2(k+1-m)/N} |u|_{k+1,K}^2,$$

where $C_1 = C_1(k,m) > 0$. Finally, note that (4.26) yields

$$(4.38) \quad |u|_{k+1,K}^2 = \sum_{|\alpha|=k+1} \int_K |D^\alpha u|^2 = \int_K 1 = \text{vol}(K),$$

and so (4.27) now follows from (4.37) and (4.38).

Hence, $\mathcal{J}_n^\circ \subset \mathcal{J}_n$ yields

$$(4.39) \quad \begin{aligned} \inf_{s \in \mathcal{S}_n} |u - s|_m^2 &\geq \sum_{K \in \mathcal{J}_n^\circ} \inf_{s \in P_k(K)} |u - s|_{m,K}^2 \\ &\geq C_1^2 \sum_{K \in \mathcal{J}_n^\circ} \text{vol}(K)^{2(k+1-m)/N + 1}. \end{aligned}$$

Since

$$(4.40) \quad \sum_{K \in \mathcal{J}_n^\circ} \text{vol}(K) = \text{vol}(\overline{\Omega_n}),$$

we may use calculus to find that

$$(4.41) \quad \sum_{K \in \mathcal{J}_n^\circ} \text{vol}(K)^{2(k+1-m)/N + 1} \geq \left[\frac{\text{vol}(\overline{\Omega_n})}{\#\mathcal{J}_n^\circ} \right]^{2(k+1-m)/N}.$$

From (4.24) and (4.25), there is an $n_0 > 0$ such that

$$(4.42) \quad \text{vol}(\overline{\Omega}_n) \geq \frac{1}{2} \text{vol}(\overline{\Omega}^0) \quad \forall n \geq n_0.$$

Hence, (4.39), (4.41), and (4.42) yield that there is a $C_2 > 0$, independent of n , for which

$$(4.43) \quad \inf_{s \in \mathfrak{S}_n} |u - s|_m \geq C_2 (\#\mathfrak{J}_n^0)^{-(k+1-m)/N} \quad \forall n \geq n_0.$$

We now claim that there is a $C_3 > 0$, independent of n , such that

$$(4.44) \quad \#\mathfrak{J}_n^0 \leq C_3 n.$$

We first consider the case $m = 0$. In this case, the functions

$$(4.45) \quad \{\chi_K : K \in \mathfrak{J}_n^0\}$$

are linearly independent elements of \mathfrak{S}_n , χ_K denoting the characteristic function of K . Since $\dim \mathfrak{S}_n = n$, we have

$$(4.46) \quad \#\mathfrak{J}_n^0 = \#\{\chi_K : K \in \mathfrak{J}_n^0\} \leq n$$

for the case $m = 0$. We now assume that $m \geq 1$. Let $\mathfrak{S}_n(\overline{\Omega}_n)$ denote the restrictions of functions in \mathfrak{S}_n to $\overline{\Omega}_n$, so that

$$(4.47) \quad \dim \mathfrak{S}_n(\overline{\Omega}_n) \leq \dim \mathfrak{S}_n = n.$$

In the case $N = 1$, we may count free parameters to see that

$$(4.48) \quad k + 1 + (\#\mathfrak{J}_n^0 - 1)(k + 1 - m) = \dim \mathfrak{S}_n(\overline{\Omega}_n),$$

so that

$$(4.49) \quad \#\mathfrak{J}_n^0 \leq \frac{n - m}{k + 1 - m}.$$

which implies (4.44) for the case $N = 1$. To establish the case for $N \geq 2$, we first claim that there is a $c_6 > 0$, independent of n , for which

$$(4.50) \quad \#\mathcal{T}_n^{\circ} \leq c_6 v(\mathcal{T}_n^{\circ}),$$

where $v(\mathcal{T}_n^{\circ})$ is the number of vertices in the triangulation \mathcal{T}_n° . Indeed, regularity of $\{\mathcal{T}_n\}_{n \geq 1}$ (and hence $\{\mathcal{T}_n^{\circ}\}_{n \geq 1}$) yields that there is a $c_6 > 0$, independent of n , such that if v is a vertex of \mathcal{T}_n° , then v can belong to at most c_6 simplices in \mathcal{T}_n° , which implies (4.50). We need only show that there is a $c_7 > 0$, independent of n , for which

$$(4.51) \quad v(\mathcal{T}_n^{\circ}) \leq c_7 \dim \mathcal{S}_n(\overline{\Omega}_n)$$

(in the case $m \geq 1$, $N \geq 2$); (4.44) then follows from (4.50), (4.51), and (4.47). Now $\mathcal{S}_n \subset C^{m-1}(\Omega)$ (see proof of Lemma 4.1). In the case $N = 2$, Theorem 1 of Ženišek [11] states that at each vertex v of \mathcal{T}_n° ,

$$(4.52) \quad \varphi \mapsto D^{\alpha} \varphi(v) \text{ for } |\alpha| \leq 2(m-1)$$

must be degrees of freedom, while the case for $N > 2$ may be reduced to the case $N = 2$ by considering restrictions of functions in $\mathcal{S}_n(\overline{\Omega}_n)$ to 2-faces of simplices $K \in \mathcal{T}_n^{\circ}$. Hence, (4.51) holds with $c_7 = \binom{N + 2(m-1)}{2(m-1)}^{-1}$, which finally completes the proof of (4.44).

As a result of (4.43) and (4.44), and $\|\cdot\|_m \geq |\cdot|_m$, we see that there is a $C_4 > 0$, independent of n , and an $n_0 > 0$,

such that

$$(4.53) \quad \inf_{s \in \mathcal{S}_n} \|u - s\|_m \geq C_4 n^{-(k+1-m)/N} \quad \forall n \geq n_0.$$

Now let $f = Lu$. Since $0 \neq u \in H^{2m+r}(\Omega) \cap H_E^m(\Omega)$, we see that $0 \neq f \in H^r(\Omega)$. Let

$$(4.54) \quad f^* := f / \|f\|_r$$

and

$$(4.55) \quad u^* := Sf^* = u / \|f\|_r.$$

Recall that there exists a finite $C_5 > 0$ such that $\|\cdot\|_m \leq C_5 \|\cdot\|_B$ on $H_E^m(\Omega)$. Since $\|f^*\|_r \leq 1$, we have

$$(4.56) \quad \begin{aligned} C_5 e(\varphi_n^*) &\geq \sup_{\|f\|_r \leq 1} \inf_{s \in \mathcal{S}_n} \|Sf - s\|_m \\ &\geq \inf_{s \in \mathcal{S}_n} \|Sf^* - s\|_m = \inf_{s \in \mathcal{S}_n} \|u^* - s\|_m \\ &= \frac{1}{\|f\|_r} \inf_{s \in \mathcal{S}_n} \|u - s\|_m \quad (\text{since } \mathcal{S}_n \text{ is a subspace}) \\ &\geq \frac{C_4}{\|f\|_r} n^{-(k+1-m)/N} \quad \forall n \geq n_0, \end{aligned}$$

which establishes (4.17) and the theorem. ■

We now ask whether the FEM is asymptotically optimal using the information n_n^* . We find that this is the case if and only if $k \geq 2m - 1 + r$ from

Theorem 4.2.

- (i) $r(n_n^*) = \Omega(n^{-(m+r)/N})$ as $n \rightarrow \infty$.
- (ii) If $\{\mathfrak{F}_n\}_{n=1}^{\infty}$ is quasi-uniform, then

$$(4.57) \quad e(\varphi_n^S) = r(n_n^*) = \Theta(n^{-(r+m)/N}) \text{ as } n \rightarrow \infty,$$

where φ_n^S is the spline algorithm using the information n_n^* .

Proof: Using (3.15), (3.22), (3.23), and (3.26), we find

$$(4.58) \quad r(n_n^*) \geq r(n) = \Theta(n^{-(r+m)/N}) \text{ as } n \rightarrow \infty,$$

establishing (i). To establish (ii), let $z \in \mathfrak{F}_0 \cap \ker n_n^*$.

Then

$$(4.59) \quad z \in \ker n_n^* \Rightarrow (z, s)_0 = 0 \quad \forall s \in \mathfrak{S}_n$$

and

$$(4.60) \quad z \in \mathfrak{F}_0 \Rightarrow z \in H^r(\Omega) \text{ and } \|z\|_r \leq 1.$$

From (2.8) and (2.9), we see that (4.59) yields

$$(4.61) \quad \begin{aligned} \|Sz\|_B^2 &= B(Sz, Sz) = (z, Sz)_0 \\ &= (z, Sz - s)_0 \quad \forall s \in \mathfrak{S}_n \\ &\leq \|z\|_r \|Sz - s\|_{-r} \quad \forall s \in \mathfrak{S}_n. \end{aligned}$$

By Theorem 4.1.1 of [1] and the equivalence of $\|\cdot\|_B$ and $\|\cdot\|_m$, there exists $s_n \in \mathfrak{S}_n$, such that

$$(4.62) \quad \|Sz - s_n\|_{-r} \leq C_1 n^{-\lambda/N} \|Sz\|_m \leq C_2 n^{-\lambda/N} \|Sz\|_B,$$

where C_1 and C_2 are positive constants independent of n , and

$$(4.63) \quad \lambda = \min(k + 1 + r, m + r) = m + r$$

by Lemma 4.1. Hence (4.61)-(4.63) yield that

$$(4.64) \quad \|Sz\|_B \leq C_2 n^{-(m+r)/N} \quad \forall z \in \mathcal{F}_0 \cap \ker \eta_n,$$

and so (3.16) yields

$$(4.65) \quad r(\eta_n) = \sup_{z \in \mathcal{F}_0 \cap \ker \eta_n} \|Sz\|_B \leq C_2 n^{-(m+r)/N}.$$

Using (4.58), (4.65), and (3.20), we find (ii). ■

Hence, the information η_n^* is (asymptotically) an n th optimal information. In the case that $k \geq 2m - 1 + r$, the FEM is (asymptotically) an optimal error algorithm using η_n^* ; when this inequality no longer holds, the FEM is no longer an asymptotically optimal error algorithm.

Remark 4.1. In Section 2, some rather stringent assumptions were made concerning the smoothness of the region and the coefficients appearing in the differential operators L, B_0, \dots, B_{m-1} . If these smoothness assumptions are violated, the shift theorem no longer holds; that is, although the second inequality in (2.10) holds for all $r \geq -m$, the first inequality may only hold for all r in some subinterval $[-m, r_0)$. Since the shift theorem no longer holds for all r , the error of the FEM is now $\Omega(n^{-(m+r_0)/N})$ as $n \rightarrow \infty$, no matter how big r is, and no matter how k is chosen. On the other hand, the proof of the error estimate of the spline algorithm does not use the shift theorem.

Hence, Theorem 4.2 holds, even if the smoothness conditions imposed in Section 2 are drastically weakened.

As an example, we consider a problem with mixed Dirichlet-Neumann boundary conditions. Let

$$(4.66) \quad \partial\Omega = \Gamma_D \cup \Gamma_N$$

be a partition of $\partial\Omega$ such that Γ_D is of positive boundary measure. Let

$$(4.67) \quad H_E^1(\Omega) := \{v \in H^1(\Omega) : v = 0 \text{ on } \Gamma_D\}.$$

For $r \geq -1$, we consider the problem of finding, for any $f \in H^r(\Omega)$, a function $u = Sf \in H_E^1(\Omega)$ satisfying

$$(4.68) \quad \int_{\Omega} \nabla u \cdot \nabla v = \int_{\Omega} f v \quad \forall v \in H_E^1(\Omega).$$

This is the weak form of the problem

$$(4.69) \quad \begin{aligned} -\Delta u &= f && \text{in } \Omega \\ u &= 0 && \text{on } \Gamma_D \\ \frac{\partial u}{\partial \nu} &= 0 && \text{on } \Gamma_N \end{aligned}$$

($\frac{\partial}{\partial \nu}$ denoting the normal derivative).

We wish to put (4.68) into the notation of Section 2.

Let χ_D and χ_N denote the characteristic functions of Γ_D and Γ_N , respectively. Define

$$(4.70) \quad \begin{aligned} Lv &:= -\Delta v \\ B_0 v &:= \chi_D v + \chi_N \frac{\partial v}{\partial \nu} \end{aligned} .$$

Then $u = Sf$ is the weak solution to

$$(4.71) \quad \begin{aligned} Lu &= f && \text{in } \Omega \\ B_0 u &= 0 && \text{on } \partial\Omega \end{aligned}$$

However, the coefficients appearing in B_0 are discontinuous, i.e., the smoothness assumptions of Section 2 are violated.

From results of Grisvard [4], the first inequality in (2.10) holds only for $r \in [-1, \frac{1}{2})$, i.e., $r_0 = \frac{1}{2}$ above. Hence, the FEM has error $O(n^{-3/(2N)})$ as $n \rightarrow \infty$, while the spline algorithm has error $O(n^{-(m+r)/N})$ as $n \rightarrow \infty$. ■

5. Local Error and Generalized Galerkin Methods

In this section, we discuss generalized Galerkin methods (i.e., with different spaces of test and trial functions) for the variational boundary-value problem. We wish to determine when such a method has local error which is bounded by a constant multiple of the optimal local error. In particular, we will exhibit a FEM which is (worst-case) asymptotically optimal, yet whose local error is arbitrarily worse than the optimal local error.

To measure the amount by which the local error of an algorithm varies from the optimal local error, Traub, Wasilkowski, and Woźniakowski (Chapter 4 of [9]) introduced the concept of "deviation." Let φ be an algorithm using n . Then the deviation $\text{dev}(\varphi)$ of φ is defined to be

$$(5.1) \quad \text{dev}(\varphi) := \sup_{f \in \mathcal{F}_0} \frac{e(\varphi, f)}{\text{rad } Uf} .$$

Clearly

$$(5.2) \quad \text{dev}(\varphi) \geq 1,$$

with

$$(5.3) \quad \text{dev}(\varphi) = 1 \quad \text{iff } \varphi \text{ is strongly optimal error.}$$

Moreover, Traub et al. showed that if φ is homogeneous, i.e.,

$$(5.4) \quad \tilde{\varphi}(\alpha y) = \alpha \varphi(y) \quad \forall \alpha \in \mathbb{R}, \quad y \in \mathbb{R}^n$$

then (in our Hilbert space setting)

$$(5.5) \quad \text{dev}(\varphi) < \infty \quad \text{iff} \quad \varphi = \varphi^S.$$

So, the spline method is the only homogeneous algorithm having finite deviation.

We wish to determine when the FEM has finite deviation. Since the FEM is linear (and hence homogeneous), the FEM has finite deviation if and only if the FEM is a spline algorithm. Hence, we wish to investigate when the FEM is a spline algorithm. In order to expedite this investigation, we now define generalized Galerkin methods (which include standard Galerkin methods and hence FEM's).

Let $\{s_i\}_{i=1}^n$ and $\{t_i\}_{i=1}^n$ each be linearly independent sets of functions in $H_E^m(\Omega)$. Let

$$(5.6) \quad \mathfrak{S} := \text{span}\{s_i\}_{i=1}^n \quad \text{and} \quad \mathfrak{T} := \text{span}\{t_i\}_{i=1}^n$$

denote the subspaces of test and trial functions (respectively).

We define the generalized Galerkin method $\varphi_{\mathfrak{S}, \mathfrak{T}}$ using \mathfrak{S} and \mathfrak{T} by

$$(5.7) \quad \varphi_{\mathfrak{S}, \mathfrak{T}}(n_{\mathfrak{S}} f) := u_{\mathfrak{S}, \mathfrak{T}},$$

where $u_{\mathfrak{S}, \mathfrak{T}} \in \mathfrak{T}$ satisfies

$$(5.8) \quad B(u_{\mathfrak{S}, \mathfrak{T}}, s) = (f, s)_0 \quad \forall s \in \mathfrak{S}$$

and $n_{\mathfrak{S}}$ is the Galerkin information (3.2) generated by \mathfrak{S} .

Remark 5.1. The (standard) Galerkin method $\varphi_{\mathfrak{S}}$ is a generalized Galerkin method with $\mathfrak{T} = \mathfrak{S}$. The FEM is a generalized Galerkin method with $\mathfrak{T} = \mathfrak{S} = \mathfrak{S}_n$, with \mathfrak{S}_n an n -dimensional subspace of $H_E^m(\Omega)$ consisting of piecewise polynomials of degree k . ■

Remark 5.2. It is more common to use Galerkin methods with different spaces of test and trial functions in the solution of variational problems associated with a bilinear form on $u \times v$, where u and v are different Hilbert spaces. See e.g. pp. 217-220 in Section 6.3 of [1] for such a method. ■

In what follows, we let $S^* : H_E^m(\Omega) \rightarrow H^r(\Omega)$ denote the Hilbert space adjoint to S , remembering that $H_E^m(\Omega)$ is a Hilbert space under the inner product given by the bilinear form B . Hence (2.9) yields

$$(5.9) \quad (g, v)_0 = B(Sg, v) = (g, S^*v)_r \quad \forall v \in H_E^m(\Omega), \quad g \in H^r(\Omega).$$

We then have

$$\text{Lemma 5.1.} \quad (\ker n)^\perp = S^*\mathfrak{S}.$$

Proof: Let $s \in \mathfrak{S}$. Then for any $h \in \ker n$,

$$(5.10) \quad (S^*s, h)_r = B(Sh, s) = (h, s)_0 = 0.$$

Hence $S^*\mathfrak{S} \subset (\ker n)^\perp$. Now S is a dense injection, so that S^* is an injection. So $\#n = n$ yields

$$(5.11) \quad \dim S^*\mathfrak{S} = \dim \mathfrak{S} = n = \dim(\ker n)^\perp,$$

which, along with $S^*\mathfrak{S} \subset (\ker n)^\perp$, yields the desired result. ■

Lemma 5.2. Given n -dimensional subspaces \mathfrak{S} and \mathfrak{T} of $H_E^m(\Omega)$, suppose that bases $\{s_i\}_{i=1}^n$ and $\{t_i\}_{i=1}^n$ of \mathfrak{S} and \mathfrak{T} , respectively, are chosen such that

$$(5.12) \quad (S^*s_j, S^*s_i)_r = \delta_{ij} \quad (1 \leq i, j \leq n)$$

and

$$(5.13) \quad B(t_j, s_i) = \delta_{ij} \quad (1 \leq i, j \leq n).$$

Then

$$(5.14) \quad \varphi^S(n_g f) = \sum_{i=1}^n (f, s_i)_0 S S^* s_i$$

and

$$(5.15) \quad \varphi_{S, \mathcal{J}}(n_g f) = \sum_{i=1}^n (f, s_i)_0 t_i.$$

Proof: Let $f_i = S^* s_i$. Then Lemma 3.1 yields $f_i \in (\ker n_g)^\perp$. In addition, (5.9) (with $g = S^* s_i$ and $v = s_j$) and (5.12) yield

$$(5.16) \quad (f_i, s_j)_0 = (S^* s_i, s_j)_0 = (S^* s_i, S^* s_j)_r = \delta_{ij}.$$

The representation formula (5.14) for the spline algorithm φ^S now follows from (5.1) of Chapter 4 of [9].

To see (5.15), write $u_{S, \mathcal{J}}$ in the form

$$(5.17) \quad u_{S, \mathcal{J}} = \sum_{j=1}^n \alpha_j t_j.$$

Then (5.8) and (5.13) yield

$$(5.18) \quad (f, s_i)_0 = B(u_{S, \mathcal{J}}, s_i) = \sum_{j=1}^n \alpha_j B(t_j, s_i) = \alpha_i,$$

establishing (3.15). ■

We now give the main result of this section, which tells us the unique choice of trial function space \mathcal{J} (corresponding to

the given test function space \mathcal{S}) for which the generalized Galerkin method is the spline method, i.e., for which the generalized Galerkin method has finite deviation.

Theorem 5.1. Let \mathcal{S} and \mathcal{T} be n -dimensional subspaces of $H_E^m(\Omega)$. Then the following are equivalent:

- (i) $\text{dev}(\varphi_{\mathcal{S}, \mathcal{T}}) < \infty$.
- (ii) $\varphi_{\mathcal{S}, \mathcal{T}} = \varphi^{\mathcal{S}}$.
- (iii) $\mathcal{T} = \mathcal{S}\mathcal{S}^*$.

Proof: Since $\varphi_{\mathcal{S}, \mathcal{T}}$ is linear and thus homogeneous, (i) and (ii) are equivalent by (5.5). We show that (ii) and (iii) are equivalent. Let \mathcal{S} and \mathcal{T} be n -dimensional subspaces of $H_E^m(\Omega)$; choose a basis $\{s_i\}_{i=1}^n$ for \mathcal{S} such that (5.12) holds.

Suppose first that (ii) holds. Choosing a basis $\{t_i\}_{i=1}^n$ for \mathcal{T} such that (5.13) holds, Lemma 5.1 yields (5.14) and (5.15). Using (5.9), we have

$$\begin{aligned} t_i &= \sum_{j=1}^n (S^* s_i, s_j)_0 t_j = \varphi_{\mathcal{S}, \mathcal{T}}(n_{\mathcal{S}} S^* s_i) \\ (5.19) \quad &= \varphi^{\mathcal{S}}(n_{\mathcal{S}} S^* s_i) = \sum_{j=1}^n (S^* s_i, s_j)_0 \mathcal{S}\mathcal{S}^* s_j = \mathcal{S}\mathcal{S}^* s_i \end{aligned}$$

for $1 \leq i \leq n$, so that (5.6) yields $\mathcal{T} = \mathcal{S}\mathcal{S}^*$. So, (ii) implies (iii).

Now suppose that (iii) holds. Let

$$(5.20) \quad t_i = \mathcal{S}\mathcal{S}^* s_i \quad (1 \leq i \leq n).$$

Then (iii) and the injectivity of $\mathcal{S}\mathcal{S}^*$ show that $\{t_i\}_{i=1}^n$ is a basis for \mathcal{T} . Using (5.9) and (5.12), we have (for

$$1 \leq i, j \leq n$$

$$(5.21) \quad B(t_j, s_i) = B(SS^* s_j, s_i) = (S^* s_j, S^* s_i)_r = \delta_{ij},$$

so that (5.13) and (5.14) hold. Using (5.13), (5.14), and (5.15), we see that (ii) holds. Thus (iii) implies (ii). \square

Hence, given any finite-dimensional subspace \mathcal{S} of $H_E^m(\Omega)$ we see how to choose the unique subspace \mathcal{T} of $H_E^m(\Omega)$ with $\dim \mathcal{T} = \dim \mathcal{S}$ such that $\varphi^{\mathcal{S}} = \varphi_{\mathcal{S}, \mathcal{T}}$. On the other hand, the most natural choice of subspace is to pick $\mathcal{T} = \mathcal{S}$, so that we get the standard Galerkin method $\varphi_{\mathcal{S}}$. When is $\text{dev}(\varphi_{\mathcal{S}})$ finite, i.e., when is $\varphi_{\mathcal{S}}$ the spline method?

Theorem 5.2. Let \mathcal{S} be an n -dimensional subspace of $H_E^m(\Omega)$. Then the following conditions are equivalent:

- (i) $\text{dev}(\varphi_{\mathcal{S}}) < \infty$.
- (ii) $\varphi_{\mathcal{S}} = \varphi^{\mathcal{S}}$.
- (iii) $\mathcal{S} = SS^* \mathcal{S}$.
- (iv) \mathcal{S} is an eigenspace of SS^* .
- (v) $\mathcal{S} = S\mathcal{T}$, where \mathcal{T} is an n -dimensional subspace of $H^r(\Omega)$ such that $\mathcal{T} = S^* \mathcal{T}$.
- (vi) $\mathcal{S} = S\mathcal{T}$, where \mathcal{T} is an n -dimensional eigenspace of $S^* S$.

Proof: From Theorem 5.1, we have (i), (ii), and (iii) are equivalent. Suppose that (iii) holds. Then $SS^* : \mathcal{S} \rightarrow \mathcal{S}$ is self-adjoint, so that \mathcal{S} (being finite-dimensional) has a basis of eigenvectors of SS^* , i.e., \mathcal{S} is an eigenspace of SS^* , i.e., (iii) implies (iv). On the other hand, an eigenspace

of an operator is always invariant under that operator, i.e., (iv) implies that $SS^*\mathfrak{S} \subset \mathfrak{S}$; the inclusion $\mathfrak{S} \subset SS^*\mathfrak{S}$ follows from the injectivity of SS^* and the finite-dimensionality of \mathfrak{S} . So (iii) and (iv) are equivalent.

Suppose that (iv) holds. Let $\mathfrak{F} = S^*\mathfrak{S}$. Using the equivalence of (iii) and (iv), we have $S^*S\mathfrak{F} = S^*SS^*\mathfrak{S} = S^*\mathfrak{S} = \mathfrak{F}$, while injectivity of S^* yields $\dim \mathfrak{F} = \dim \mathfrak{S} = n$. Moreover, $S\mathfrak{F} = SS^*\mathfrak{S} = \mathfrak{S}$. So (iv) implies (v). If (v) holds, then $\mathfrak{S} = S\mathfrak{F} = SS^*S\mathfrak{F} = SS^*\mathfrak{S}$, so (v) implies (iii), which in turn yields (iv). So (iv) and (v) are equivalent.

Finally, $\mathfrak{F} = S^*S\mathfrak{F}$ if and only if \mathfrak{F} is an eigenspace of S^*S , the argument being similar to that in the preceding paragraph. Hence (v) and (vi) are equivalent. ■

We now consider two examples for which one of the conditions in Theorem 5.2 holds, so that the Galerkin method and the spline method are one and the same.

Example 5.1. Let $r = -m$. Then S is the Riesz map, which is an isometric isomorphism of $H^{-m}(\Omega)$ (under the norm $\|S \cdot\|_B$, which is equivalent to $\|\cdot\|_{-m}$) onto $H_E^m(\Omega)$ (under the norm $\|\cdot\|_B$), see Section 4.4 of [7]. Hence $SS^* = I$, the identity map on $H_E^m(\Omega)$, and so $\mathfrak{S} = SS^*\mathfrak{S}$ for any subspace \mathfrak{S} of $H_E^m(\Omega)$. So when $r = -m$, the standard Galerkin method is the spline algorithm, no matter what the choice of \mathfrak{S} . Of course in this case, (3.26) shows that $\lim_{n \rightarrow \infty} r(n) \neq 0$, i.e., there is no convergent sequence of algorithms, each of which uses finite information; see also Corollary 5.1 in Chapter 2 of [9]. ■

Example 5.2. Let $r > -m$. Then $H^r(\Omega)$ has a complete basis of eigenvectors for S^*S , i.e., there exists a basis $\{e_i\}_{i=1}^{\infty}$ for $H^r(\Omega)$ and real numbers $\lambda_1 \geq \lambda_2 \geq \dots > 0$ with $\lim_{i \rightarrow \infty} \lambda_i = 0$, such that for any positive integers i and j ,

$$(5.22) \quad S^* S e_i = \lambda_i e_i$$

and

$$(5.23) \quad (e_j, e_i)_r = \delta_{ij}.$$

Define the information $n_n : H^r(\Omega) \rightarrow \mathbb{R}^n$ by

$$(5.24) \quad n_n f := [(f, e_1)_r \dots (f, e_n)_r]^T.$$

Then $\#n_n = n$. Moreover, n_n is the n th optimal information, and

$$(5.25) \quad r(n_n) = r(n) = \sqrt{\lambda_{n+1}} = \Theta(n^{-(r+m)/N}) \quad \text{as } n \rightarrow \infty$$

(Theorem 5.3 of Chapter 2 of [9]). Letting

$$(5.26) \quad \mathfrak{F} = \text{span}\{e_1, \dots, e_n\},$$

we see that (vi) of Theorem 5.2 holds. Setting

$$(5.27) \quad \mathfrak{S} = \text{span}\{s_1, \dots, s_n\}, \quad s_i := \frac{1}{\lambda_i} S e_i \quad (1 \leq i \leq n),$$

we find $n_n = n_{\mathfrak{S}}$, since (4.27), (4.32), and (4.5) yield

$$(5.28) \quad (f, e_i)_r = \left(f, \frac{S^* S e_i}{\lambda_i} \right)_r = (f, S^* s_i)_r = (f, s_i)_0 \quad (1 \leq i \leq n).$$

Hence the spline method and the standard Galerkin method coincide

for the n th optimal information n_n . Since

$$(5.29) \quad (S^* s_j, S^* s_i)_r = \frac{1}{\lambda_i \lambda_j} (S^* S e_j, S^* S e_i)_r = (e_j, e_i)_r = \delta_{ij},$$

we see that the formula for the spline method and standard Galerkin method is given by

$$(5.30) \quad \varphi^S(n_n f) = \varphi_S(n_n f) = \sum_{i=1}^n (f, s_i)_0 S S^* s_i = \sum_{i=1}^n \lambda_i (f, e_i)_r S e_i$$

in this case. ■

We now turn to an example $S \neq S S^*$. This example is of particular interest because it gives an instance of an FEM which has optimal worst-case error (to within a constant, independent of n), but has infinite deviation.

Example 5.3. We consider the L_2 -approximation problem for H^1 -functions on the interval $(0,1)$. Choose $N = 1$, $m = 0$, $r = 1$, and let $S : H^1(0,1) \rightarrow L_2(0,1)$ be the canonical injection

$$(5.31) \quad Su := u \quad \forall u \in H^1(0,1).$$

The variational form of the problem is to define

$B : L_2(0,1) \times L_2(0,1) \rightarrow \mathbb{R}$ by

$$(5.32) \quad B(u,v) := \int_0^1 uv \quad \forall u,v \in L_2(0,1).$$

Then for any $f \in H^1(0,1)$, we wish to find $u = Sf \in L_2(0,1)$ such that

$$(5.33) \quad B(u,v) = \int_0^1 fv \quad \forall v \in L_2(0,1).$$

(Of course, $u = f$.)

We let \mathcal{S}_n be an n -dimensional subspace of $L_2(0,1)$ consisting of piecewise constants, so that $k = 0$. Let

$$(5.34) \quad 0 = x_0 < x_1 < \dots < x_{n-1} < x_n = 1$$

be a partition of $(0,1)$. Then \mathcal{S}_n is the span of the functions s_1, \dots, s_n where

$$(5.35) \quad s_i(x) = \delta_{ij} \quad x_{j-1} \leq x \leq x_j \quad (1 \leq j \leq n, 1 \leq i \leq n).$$

Using an integration by parts, one can show that for any $s \in L_2(\Omega)$, $w := SS^*s$ is the (weak) solution to

$$(5.36) \quad \begin{aligned} -w'' + w &= s && \text{in } (0,1) \\ w'(0) &= w'(1) = 0 \end{aligned}$$

so that

$$(5.37) \quad w(x) = \frac{\int_0^1 s(\xi) \cosh(1 - \xi) d\xi}{\sinh 1} \cosh x - \int_0^x s(\xi) \sinh(x - \xi) d\xi.$$

Hence $SS^*\mathcal{S}_n$ is the span of $\{w_1, \dots, w_n\}$, where

$$(5.38) \quad \begin{aligned} w_i(x) &= \frac{\kappa_i}{\sinh 1} \cosh x - \int_0^1 s_i(\xi) \sinh(x - \xi) d\xi \\ &= \frac{\kappa_i}{\sinh 1} \cosh x - [\cosh(x - x_{i-1}) - \cosh(x - x_i)] \end{aligned}$$

and

$$(5.39) \quad \kappa_i = \int_0^1 s_i(\xi) \cosh(1 - \xi) d\xi = \sinh(1 - x_{i-1}) - \sinh(1 - x_i).$$

Since none of the w_i is piecewise constant on $(0,1)$, we have $w_i \notin \mathcal{S}_n$, so that $SS^*\mathcal{S}_n \neq \mathcal{S}_n$.

Hence, the FEM is not the spline method in this case, and thus has infinite deviation. This is interesting for the following version. Suppose that $\{\mathcal{S}_n\}_{n=1}^{\infty}$ is quasi-uniform. Then the FEM has worst-case error $\Theta(n^{-1})$, and is (to within a constant, independent of n) an optimal-error algorithm. Hence, we have an example of an almost optimal-error algorithm that has infinite deviation, i.e., is not strongly optimal error. ■

Examples 5.1 and 5.3 suggest the following

Conjecture 5.1. Let $r > -m$ and let \mathcal{S}_n be a finite-element subspace of $H_E^m(\Omega)$. Then the FEM using \mathcal{S}_n has infinite deviation. ■

From the results of Section 4, it is clear that Conjecture 5.1 holds when $k < 2m - 1 + r$. Hence, it remains only to prove the conjecture for the case $k \geq 2m - 1 + r$.

6. Complexity Analysis

In this section, we discuss the complexity of finding ϵ -approximations to the solution of the variational boundary-value problem, as well as the penalty for using the FEM when $k < 2m - 1 + r$.

Let $\epsilon > 0$. An algorithm φ using n furnishes an ϵ -approximation to the problem if

$$(6.1) \quad e(\varphi) \leq \epsilon.$$

The complexity $\text{comp}(\varphi)$ of an algorithm φ using n is defined in the model of computation discussed in Chapter 5 of [9].

(Informally, we assume that linear functionals can be evaluated in finite time and that the cost of an arithmetic operation is unity.) It then turns out that for any algorithm φ using n of cardinality n ,

$$(6.2) \quad \text{comp}(\varphi) \geq nc_1 + n - 1,$$

c_1 being the complexity of evaluating a linear functional, while if φ is a linear function of the information used, then

$$(6.3) \quad \text{comp}(\varphi) \leq nc_1 + 2n - 1.$$

(See Section 2, Chapter 5 of [9] for further details.) We then define, for $\epsilon > 0$, the ϵ -complexity $\text{COMP}(\epsilon)$ of the problem to be

$$(6.4) \quad \text{COMP}(\epsilon) := \inf\{\text{comp}(\varphi) : e(\varphi) \leq \epsilon\},$$

the infimum being taken over all such φ using information of finite cardinality.

Remark 6.1. Note that we distinguish between algorithmic complexity and problem complexity. For an algorithm φ , $\text{comp}(\varphi)$ denotes the complexity of the algorithm φ , while for $\epsilon > 0$, $\text{COMP}(\epsilon)$ denotes the (minimal) complexity of finding an ϵ -approximation. To tie these two concepts together, let $\epsilon > 0$. Suppose that φ^* is an algorithm with

$$(6.5) \quad e(\varphi^*) \leq \epsilon$$

and such that for any other algorithm φ ,

$$(6.6) \quad e(\varphi) \leq \epsilon \Rightarrow \text{comp}(\varphi) \geq \text{comp}(\varphi^*).$$

Then φ^* is an optimal complexity algorithm for finding an ϵ -approximation, and

$$(6.7) \quad \text{COMP}(\epsilon) = \text{comp}(\varphi^*). \quad \blacksquare$$

Let $\{\mathcal{S}_n\}_{n=1}^{\infty}$ be a quasi-uniform family of finite element subspaces of $H_E^m(\Omega)$ consisting of piecewise polynomials of degree k . Let φ_n^* be the FEM based on the space \mathcal{S}_n ; that is, for $f \in \mathcal{F}_0$, we let $u_n^* \in \mathcal{S}_n$ satisfy

$$(6.8) \quad B(u_n^*, s) = (f, s)_0 \quad \forall s \in \mathcal{S}_n,$$

and then set

$$(6.9) \quad \varphi_n^*(h_n^* f) := u_n^*.$$

We wish to measure the algorithmic complexity of using the FEM to find ϵ -approximations, i.e.,

$$(6.10) \quad \text{FEM}(\epsilon) := \inf\{\text{comp}(\varphi_n^*) : e(\varphi_n^*) \leq \epsilon\},$$

and compare $\text{FEM}(\epsilon)$,

$$(6.11) \quad \text{SPLINE}(\epsilon) := \inf\{\text{comp}(\varphi_n^S) : e(\varphi_n^S) \leq \epsilon\}$$

(φ_n^S being the spline algorithm using the Galerkin information n_n^* generated by S_n), and $\text{COMP}(\epsilon)$.

Using (6.2), (6.4), and the results in Section 4, we find

$$(6.12) \quad \text{FEM}(\epsilon) = \Theta(\epsilon^{-N/\mu}) \text{ as } \epsilon \rightarrow 0,$$

where

$$(6.13) \quad \mu = \min(k + 1 - m, m + r),$$

while

$$(6.14) \quad \text{SPLINE}(\epsilon) = \Theta(\epsilon^{-N/(m+r)}) \text{ as } \epsilon \rightarrow 0$$

and

$$(6.15) \quad \text{COMP}(\epsilon) = \Theta(\epsilon^{-N/(m+r)}) \text{ as } \epsilon \rightarrow 0.$$

This yields

Theorem 6.1.

(i) The spline algorithm is asymptotically optimal, i.e.,

$$\text{SPLINE}(\epsilon) = \Theta(\text{COMP}(\epsilon)) = \Theta(\epsilon^{-N/(m+r)}) \text{ as } \epsilon \rightarrow 0.$$

(ii) If $k \geq 2m - 1 + r$, the FEM is asymptotically optimal, i.e.,

$$\text{FEM}(\epsilon) = \Theta(\text{COMP}(\epsilon)) = \Theta(\epsilon^{-N/(m+r)}), \text{ as } \epsilon \rightarrow 0.$$

(iii) If $k < 2m - 1 + r$, then

$$\frac{\text{FEM}(\epsilon)}{\text{COMP}(\epsilon)} = \Theta\left(\frac{\text{FEM}(\epsilon)}{\text{SPLINE}(\epsilon)}\right) = \Theta\left(\left(\frac{1}{\epsilon}\right)^{\lambda N}\right) \text{ as } \epsilon \rightarrow 0,$$

where

$$\lambda = \frac{1}{k+1-m} - \frac{1}{m+r} > 0,$$

so that

$$(6.16) \quad \lim_{\epsilon \rightarrow 0} \frac{\text{FEM}(\epsilon)}{\text{COMP}(\epsilon)} = \lim_{\epsilon \rightarrow 0} \frac{\text{FEM}(\epsilon)}{\text{SPLINE}(\epsilon)} = +\infty. \quad \blacksquare$$

Thus when k is too small for a given value of r , the asymptotic penalty for using the FEM instead of the spline algorithm is infinite. Clearly (6.16) tells us that there exists $\epsilon_0 > 0$ for which

$$(6.17) \quad \text{SPLINE}(\epsilon) < \text{FEM}(\epsilon) \text{ for } 0 < \epsilon < \epsilon_0.$$

What is the value of ϵ_0 ? If ϵ_0 is unreasonably small, it may turn out that it is more reasonable to use the FEM for "practical" values of ϵ . We determine the value of ϵ_0 for a model problem in

Example 6.1. Let $N = 1$, $\Omega = (0, \pi)$, $m = 1$, $r = 1$, $H_E^1(\Omega) = H_0^1(0, \pi)$, and consider the bilinear form $B : H_0^1(0, \pi) \times H_0^1(0, \pi) \rightarrow \mathbb{R}$ defined by

$$(6.18) \quad B(v, w) := \int_0^\pi v' w' \quad \forall v, w \in H_0^1(0, \pi).$$

Hence for $f \in H^1(0, \pi)$, $u = Sf$ is the variational solution to the problem

$$(6.19) \quad \begin{aligned} -u''(x) &= f(x) & 0 < x < \pi \\ u(0) &= u(\pi) = 0 \end{aligned}$$

We choose \mathcal{S}_n to be a subspace of $H_n^1(0, \pi)$ consisting of piecewise linear polynomials with nodes at $x_j = \frac{j\pi}{n+1}$ ($0 \leq j \leq n+1$). Hence $k = 1$; moreover, since any function in \mathcal{S}_n must vanish at the endpoints of $[0, \pi]$, we see that $\dim \mathcal{S}_n = n$.

We first give a lower bound on $e(\varphi_n^*)$, φ_n^* being the n th FEM. Let

$$(6.20) \quad f(x) := \frac{1}{\sqrt{\pi}}.$$

Then

$$(6.21) \quad \|f\|_1 = 1$$

and $Sf = u$, where

$$(6.22) \quad u(x) := \frac{1}{2} x \left(\sqrt{\pi} - \frac{x}{\sqrt{\pi}} \right).$$

Let \tilde{u}_n be the \mathcal{S}_n -interpolate of u , i.e., \tilde{u}_n is the unique function in \mathcal{S}_n for which

$$(6.23) \quad \tilde{u}_n(x_j) = u(x_j) \quad (1 \leq j \leq n).$$

Then using Newton's interpolation formula on each subinterval $[x_j, x_{j+1}]$ along with the fact that u'' is a constant, one can show that

$$(6.24) \quad \inf_{s \in \mathcal{S}_n} \|u - s\|_1 = \|u - \tilde{u}_n\|_1 = \frac{\tau}{\sqrt{12} (n+1)}.$$

Using (3.7), (6.21), (6.24), and $u = Sf$, we conclude that

$$(6.25) \quad e(\varphi_n^*) \geq \frac{\tau}{\sqrt{12} (n+1)}.$$

Now we can find a lower bound on $FEM(\epsilon)$. Let $e(\varphi_n^*) \leq \epsilon$.

Then (6.25) yields

$$(6.26) \quad n \geq \frac{\pi}{\sqrt{12}} \epsilon^{-1} - 1,$$

which, along with (6.2), gives the lower bound

$$(6.27) \quad FEM(\epsilon) \geq (c_1 + 1) \left(\frac{\pi}{\sqrt{12}} \epsilon^{-1} - 1 \right) - 1.$$

Next, we wish to give an upper bound on $e(\varphi_n^S)$, where φ_n^S is the spline algorithm using n_n^* . Since $e(\varphi_n^S) = r(n_n^*)$, it suffices to compute the radius of information. Let $z \in \ker n_n^* \cap BH^1$. Let P_n denote the orthogonal projector of $L_2(0, \pi)$ onto \mathcal{S}_n . Using (2.8), (2.9), the fact that $z \in \ker n_n^*$, and properties of the orthogonal projector, we find

$$(6.28) \quad \begin{aligned} \|Sz\|_B^2 &= B(Sz, Sz) = (z, Sz)_0 = (z, Sz - P_n Sz)_0 \\ &= (z - \tilde{z}_n, Sz - P_n Sz)_0 \\ &\leq \|z - \tilde{z}_n\|_0 \|Sz - P_n Sz\|_0 \\ &\leq \|z - \tilde{z}_n\|_0 \|Sz - (\tilde{Sz})_n\|_0 \end{aligned}$$

(where for $v \in H_0^1(0, \pi)$, \tilde{v}_n is the \mathcal{S}_n -interpolate of v as given by (6.23)). Since for any $v \in H_0^1(0, \pi)$, Theorem 2.4 of [7] states that

$$(6.29) \quad \|v - \tilde{v}_n\|_0 \leq \frac{1}{n+1} |v|_1,$$

(6.28) becomes

$$(6.30) \quad \|Sz\|_B^2 \leq \left(\frac{1}{n+1}\right)^2 |z|_1 |Sz|_1 \leq \left(\frac{1}{n+1}\right)^2 \|Sz\|_B,$$

where we have used $z \in BH^1$ (so that $|z|_1 \leq 1$) and $\|\cdot\|_B = |\cdot|_1$.

Hence

$$(6.31) \quad \|Sz\|_B \leq \left(\frac{1}{n+1}\right)^2 \quad \forall z \in \ker n_n^* \cap BH^1,$$

so that

$$(6.32) \quad e(\varphi_n^S) = r(n_n^*) \leq \left(\frac{1}{n+1}\right)^2$$

by (3.16).

Using (6.32), we now find an upper bound on $SPLINE(\epsilon)$.

Let

$$(6.33) \quad n \leq \epsilon^{-1/2} - 1.$$

Then (6.32) yields that $e(\varphi_n^S) \leq \epsilon$. From (6.3), we find that

$$(6.34) \quad SPLINE(\epsilon) \leq (c_1 + 2)(\epsilon^{-1/2} - 1) - 1.$$

We now wish to find $\epsilon_0 = \epsilon_0(c_1)$ such that (6.17) holds. From (6.17), (6.27), and (6.34), we see that we may choose ϵ_0 to be the smallest positive solution of

$$(6.35) \quad (c_1 + 1) \left(\frac{\pi}{\sqrt{12}} \epsilon_0^{-1} - 1 \right) = (c_1 + 2)(\epsilon_0^{-1/2} - 1).$$

Some algebra yields

$$(6.36) \quad \epsilon_0(c_1) = \left[\frac{1}{2} c_1 + 1 - \sqrt{\left(\frac{1}{2} c_1 + 1\right)^2 - \frac{\pi}{\sqrt{12}} (c_1 + 1)} \right]^2 .$$

Some elementary calculus tells us that ϵ_0 is an increasing function. Since $c_1 \geq 0$, we thus have

$$(6.37) \quad \epsilon_0(c_1) \leq \epsilon_0(0) = \left(1 - \sqrt{1 - \frac{\pi}{\sqrt{12}}} \right)^2 \doteq 0.482853424 .$$

Thus (6.17) holds for all ϵ roughly less than one-half.

On the other hand, if we are willing to assume that evaluating a linear functional is at least as hard as an arithmetic operation, we have $c_1 \geq 1$ and so

$$(6.38) \quad \epsilon_0(c_1) \leq \epsilon_0(1) = \left(\frac{3}{2} - \sqrt{\frac{9}{4} - \frac{\pi}{\sqrt{3}}} \right)^2 \doteq 0.7048360247 .$$

Of course, it is reasonable to suspect that $c_1 \gg 1$ (see e.g. pg. 85 of [9]). One may check that

$$(6.39) \quad \lim_{c_1 \rightarrow \infty} \epsilon_0(c_1) = \frac{\pi}{12} \doteq 0.8224670334 ,$$

giving an estimate of $\epsilon_0(c_1)$ for large values of c_1 . ■

Based on this example, it seems reasonable to conjecture that for any regularly-elliptic boundary-value problem, (6.17) will hold, where ϵ_0 is sufficiently large to be of interest. We suspect that such a result will be difficult to establish. There are two reasons for this. The first reason is that "sufficiently large" may be a subjective criterion. That is,

$\epsilon_0 = 10^{-1000}$ is obviously too small to be of practical interest, while $\epsilon_0 = 10^{-1}$ is not so absurdly small; where does one draw the line separating the reasonable values of ϵ_0 from the unreasonable values?

The second reason is perhaps more crucial. In order to determine ϵ_0 , we have to change the order-of-magnitude estimates in Theorems 4.1 and 4.2 to sharp bounds involving constants whose values are explicitly known. Since the values of the constants appearing in these theorems are not explicitly known in general, we suspect that this task will be very difficult for a general problem, making it very difficult to determine, for a general problem, a value of ϵ_0 such that (6.17) will hold.

7. Some Practical Considerations

In this section, we discuss the implementation of the spline algorithm. In particular, we ask when it is of practical use.

The main problem is that the spline algorithm is hard to implement. This is mainly due to the presence of the adjoint S^* to the solution operator S . As we saw in Example 5.3, even for very simple solution operators S and subspaces \mathcal{S}_n of test functions, the space $SS^*\mathcal{S}_n$ of trial functions can be very complicated.

To a certain extent, this problem may be solved by preconditioning. Generally speaking, the problem S , the class \mathcal{F}_0 of problem elements, and the family $\{\mathcal{S}_n\}_{n=1}^{\infty}$ of subspaces of test functions will be fixed. Suppose, for a given $\epsilon > 0$, we wish to compute ϵ -approximations to Sf for many $f \in \mathcal{F}_0$. Then we may determine a fixed cardinality n , depending on ϵ , such that $e(\varphi_n^S) \leq \epsilon$, i.e., the spline method φ_n^S yields ϵ -approximations for any $f \in \mathcal{F}_0$. Now, we may precondition: instead of finding $t_i = SS^*s_i$, we find approximations \tilde{t}_i to t_i for $1 \leq i \leq n$. (This may perhaps be done via an FEM.) Moreover, we may use standard techniques (e.g., the Q - R method) to biorthonormalize $\{s_i\}_{i=1}^n$ and $\{\tilde{t}_i\}_{i=1}^n$. Although this may be a lot of work, it is independent of the choice of f . Hence, if we wish to compute ϵ -approximations to Sf for many different $f \in \mathcal{F}_0$, this may be a feasible technique. (But note that since this is a linear method which does not exactly coincide with the spline algorithm, its deviation is infinite, no matter how close \tilde{t}_i and t_i are!)

On the other hand, suppose we wish to compute a sequence of approximations to Sf for a fixed $f \in \mathcal{F}_0$. In this case, the preconditioning will be prohibitively expensive, because as ϵ changes, the value of n such that $e(\varphi_n^*) \leq \epsilon$ changes, which implies that the algorithm φ_n^S changes. Since SS^*s_i ($1 \leq i \leq n$) cannot be explicitly computed for general S , it appears that the spline method will not be practical in this case. Using Theorem 4.1, it appears that the best advice is to use an FEM of sufficiently high degree, unless the problems involved in implementing such a method are so great (or ϵ is so large) that one doesn't mind the penalty of $\Theta(\epsilon^{-\lambda N})$ as $\epsilon \rightarrow 0$, where $\lambda = (k + 1 - m)^{-1} - (m + r)^{-1}$, which will result from using an FEM of degree $k < 2m - 1 + r$.

In summary, we see that for the case of solving problems $Lu = f$ with many different f to within a fixed error criterion ϵ , the spline algorithm may be of practical interest. On the other hand, we do not currently know how to efficiently implement the spline algorithm when solving a single problem with greater and greater accuracy.

Acknowledgments

I would like to thank Professors A. K. Aziz and T. I. Seidman (University of Maryland Baltimore County) and J. F. Traub (Columbia University) for their comments and suggestions.

References

- [1] I. Babuska and A. K. Aziz, "Survey lectures on the mathematical foundations of the finite element method," in The Mathematical Foundations of the Finite Element Method with Applications to Partial Differential Equations, (A. K. Aziz, ed.), Academic Press, New York, 1972, pp. 3-359.
- [2] P. G. Ciarlet, The Finite Element Method for Elliptic Problems, North-Holland, Amsterdam, 1978.
- [3] P. G. Ciarlet and P. A. Raviart, "Interpolation theory over curved elements," Comput. Methods Appl. Mech. Engrg., v. 1, 1972, pp. 217-249.
- [4] P. Grisvard, "Behavior of the solutions of an elliptic boundary-value problem in a polygonal or polyhedral domain," in Numerical Solutions of Partial Differential Equations -- III (SYNSPADE 1975), (Bert Hubbard, ed.), Academic Press, New York, 1976, pp. 207-274.
- [5] D. E. Knuth, "Big Omicron and Big Omega and Big Theta," SIGACT News, ACM, April, 1976.
- [6] J. T. Oden and J. N. Reddy, An Introduction to the Mathematical Theory of Finite Elements, Wiley-Interscience, New York, 1976.
- [7] M. Schultz, Spline Analysis, Prentice-Hall, Englewood Cliffs, 1973.
- [8] G. Strang and G. Fix, "A Fourier analysis of the finite element variational method," in Constructive Aspects of Functional Analysis, Part II, C.I.M.E., Rome, 1973.
- [9] J. F. Traub and H. Woźniakowski, A General Theory of Optimal Algorithms, Academic Press, New York, 1980.
- [10] A. G. Werschulz, "Optimal error properties of finite element methods for second order elliptic Dirichlet problems," Math. Comp., v. 38, No. 158, April, 1982, pp. 401-413.
- [11] A. Ženišek, "Hermite interpolation on simplexes in the finite element method," in Proceedings EquaDiff 3, J. E. Purkyně University, Brno, pp. 271-277.