

# On the condition number of certain Rayleigh-Ritz-Galerkin matrices

Bernard J. Omodei

Martin H. Schultz [*Bull. Amer. Math. Soc.* 76 (1970), 840-844] has investigated the spectral condition number of the Rayleigh-Ritz-Galerkin matrices that arise when normalized  $B$ -spline coordinate functions are used to approximate the solution of a class of linear, self-adjoint, elliptic boundary value problems in one dimension. This paper shows how results analogous to those of Schultz [*op. cit.*] can be established under weaker assumptions. We also extend the results to boundary value problems in higher dimensions.

We consider the following class of linear, self-adjoint, two-point boundary value problems:

$$(1) \quad L[u(x)] \equiv \sum_{j=0}^n (-1)^j D^j \left[ p_j(x) D^j u(x) \right] = f(x) ,$$

$$0 < x < 1 , \quad f \in L^2[0, 1] , \quad n \geq 1 ,$$

with homogeneous Dirichlet boundary conditions

$$(2) \quad D^k u(0) = D^k u(1) = 0 , \quad 0 \leq k \leq n - 1 .$$

Assume that  $p_j(x)$  ,  $0 \leq j \leq n$  , are real-valued bounded measurable functions on  $[0, 1]$  .

Let  $W_0^{n,2}[0, 1]$  denote the completion of the set of all  $C^\infty[0, 1]$  functions having compact support in  $(0, 1)$  , with respect to the Sobolev

---

Received 17 June 1976.

norm

$$\|w\|_{W^{n,2}} \equiv \left\{ \sum_{j=0}^n \int_0^1 [D^j w(x)]^2 dx \right\}^{\frac{1}{2}}.$$

We assume that there exists a positive constant  $K$  such that for all  $w \in W_0^{n,2}[0, 1]$ ,

$$(3) \quad K \|w\|_{L^2}^2 \leq \int_0^1 \left\{ \sum_{j=0}^n p_j(x) [D^j w(x)]^2 \right\} dx.$$

This assumption corresponds to the assumption that the differential operator  $L$  is positive definite. Schultz [7] made the stronger assumption that, for all  $w \in W_0^{n,2}[0, 1]$ ,

$$K \|w\|_{W^{n,2}}^2 \leq \int_0^1 \left\{ \sum_{j=0}^n p_j(x) [D^j w(x)]^2 \right\} dx.$$

It can be shown that the problem (1)-(2) has a unique generalized solution and that the Rayleigh-Ritz method is applicable; see Omodei [6].

Let  $\{\phi_i(x)\}_{i=1}^m$  be  $m$  given linearly independent coordinate functions such that  $\phi_i \in W_0^{n,2}[0, 1]$  for all  $1 \leq i \leq m$ . Let  $S_m$  denote the approximating subspace spanned by  $\{\phi_i\}_{i=1}^m$ . We claim, without giving the derivation, that the Rayleigh-Ritz-Galerkin matrix  $R \equiv (r_{ik})$  for the problem (1)-(2) is given by

$$(4) \quad r_{ik} = \int_0^1 \left\{ \sum_{j=0}^n p_j(x) D^j \phi_k(x) D^j \phi_i(x) \right\} dx, \quad 1 \leq i, k \leq m.$$

We now introduce normalized  $B$ -spline coordinate functions. Following the construction of de Boor [1], for a positive integer  $d$ , the finite set of real numbers

$$\pi : 0 = x_0 < x_1 \leq x_2 \leq \dots \leq x_N < x_{N+1} = 1$$

is said to be a  $(d+1)$ -extended partition of  $[0, 1]$ , if and only if  $x_k < x_{k+d}$  for all  $0 \leq k \leq N - d + 1$ ; that is, if  $f_k$  denotes the

multiplicity of the knot  $x_k$  in  $\pi$ , then  $f_k \leq d$  for all  $1 \leq k \leq N$ .

Let  $I \equiv \{0 \leq k \leq N \mid x_k < x_{k+1}\}$ , and define

$$(5) \quad \Delta \equiv \max_{k \in I} (x_{k+1} - x_k) \quad \text{and} \quad \delta \equiv \min_{k \in I} (x_{k+1} - x_k).$$

Let  $Sp_0(d, \pi)$  denote the *extended spline space* of all extended splines of degree  $d$  on  $\pi$  satisfying the boundary conditions (2); that is,  $Sp_0(d, \pi)$  consists of those real-valued functions on  $[0, 1]$  which satisfy the boundary conditions (2), reduce to a polynomial of degree less than or equal to  $d$  on  $[x_k, x_{k+1}]$  for all  $k \in I$ , and have  $d - f_k$  continuous derivatives in a neighbourhood of  $x_k$  for all  $1 \leq k \leq N$ .

Assuming that  $n \leq d$ , we add  $2(d-n)$  extra knots to  $\pi$  to form the partition

$$\tilde{\pi} : x_{-d+n} = \dots = x_{-1} = x_0 < x_1 \leq \dots \leq x_N < x_{N+1} = x_{N+2} = \dots = x_{N+d+1-n}.$$

We now define the classical *B-splines* for the partition  $\tilde{\pi}$  (see [4]):

$$M_k(x) \equiv (d+1)g(x_k, x_{k+1}, \dots, x_{k+d+1}; x), \quad -d + n \leq k \leq N - n,$$

is  $(d+1)$  times the  $(d+1)$ -th divided difference in  $y$  of the function

$$g(y; x) \equiv (y-x)_+^d \quad \text{based on the points } x_k, x_{k+1}, \dots, x_{k+d+1}.$$

The *normalized B-splines* are defined by

$$(6) \quad \psi_k(x) \equiv \frac{x_{k+d+1} - x_k}{d+1} M_k(x), \quad -d + n \leq k \leq N - n.$$

It can be shown that  $\{\psi_k(x)\}_{k=-d+n}^{N-n}$  form a basis for  $Sp_0(d, \pi)$  (see [4]).

The following lemma is a simple consequence of a theorem in [2].

**LEMMA 1.** *For an arbitrary  $(d+1)$ -extended partition  $\pi$ , there exists a positive constant  $D$  depending on  $d$  but not on  $\pi$  such that*

$$(7) \quad \left\| \sum_{k=-d+n}^{N-n} a_{k+d+1-n} \left( \frac{d+1}{x_{k+d+1} - x_k} \right)^{\frac{1}{2}} \psi_k \right\|_{L^2} \geq D \|a\|_2$$

for all  $a \in R^{N+d+1-2n}$  where  $\|a\|_2 \equiv \left( \sum_{i=1}^{N+d+1-2n} a_i^2 \right)^{\frac{1}{2}}$ .

We consider the case where the approximating subspace  $S_m \equiv Sp_0(d, \pi)$ ,  $m = N + d + 1 - 2n$ , and the coordinate functions  $\phi_i(x) \equiv \psi_{i+n-d-1}(x)$ ,  $i = 1, 2, \dots, m$ . Assume that  $f_k \leq d + 1 - n$  for all  $i \leq k \leq N$  to ensure that  $Sp_0(d, \pi) \subset W_0^{n,2}[0, 1]$ . The spectral condition number of the Rayleigh-Ritz-Galerkin matrix  $R$  is defined by

$$\kappa(R) \equiv \|R\|_2 \|R^{-1}\|_2 \text{ where } \|R\|_2 \equiv \sup_{a \in R^m} \|Ra\|_2 / \|a\|_2.$$

Using (3), it can easily be shown that  $R$  is positive definite and symmetric, and hence  $\kappa(R) = \lambda^{-1}\Lambda$  where  $\lambda$  and  $\Lambda$  are the minimum and maximum eigenvalues, respectively, of  $R$ . The following theorem is analogous to that of Schultz [7].

**THEOREM 1.** *If (3) holds and  $\pi$  is an arbitrary  $(d+1)$ -extended partition of  $[0, 1]$  such that  $f_k \leq d + 1 - n$  for all  $1 \leq k \leq N$ , then there exists a positive constant  $C$  depending on  $d$  but not on  $\pi$  such that*

$$(8) \quad \kappa(R) \leq C(\Delta/\delta)\delta^{-2n}.$$

*Proof.* From (4) and (3), we obtain for all  $a \in R^m$ ,

$$a^T Ra = \int_0^1 \left\{ \sum_{j=0}^n p_j(x) \left[ D^j \sum_{i=1}^m a_i \psi_{i+n-d-1}(x) \right]^2 \right\} dx \geq K \left\| \sum_{i=1}^m a_i \psi_{i+n-d-1} \right\|_{L^2}^2$$

which, by Lemma 1, yields

$$\begin{aligned} a^T Ra &\geq KD^2 \sum_{i=1}^m a_i^2 \frac{(x_{i+n} - x_{i+n-d-1})}{d+1} \\ &\geq KD^2(d+1)^{-1}\delta \|a\|_2^2, \end{aligned}$$

and thus

$$(9) \quad \lambda \geq KD^2(d+1)^{-1}\delta.$$

Conversely, since  $p_j(x)$ ,  $0 \leq j \leq n$ , are bounded on  $[0, 1]$ , there exists a positive constant  $P$  such that, for all  $a \in R^m$ ,

$$\begin{aligned} a^T R a &\leq P \sum_{j=0}^n \int_0^1 \left[ \sum_{i=1}^m a_i D^j \psi_{i+n-d-1}(x) \right]^2 dx \\ &\leq P \sum_{j=0}^n (2d+1) \sum_{i=1}^m a_i^2 \int_0^1 \left[ D^j \psi_{i+n-d-1}(x) \right]^2 dx, \end{aligned}$$

since  $\psi_{i+n-d-1}(x)$ ,  $1 \leq i \leq m$ , has support  $[x_{i+n-d-1}, x_{i+n}]$ . Thus

$$a^T R a \leq P(2d+1) \sum_{i=1}^m a_i^2 (x_{i+n} - x_{i+n-d-1}) \sum_{j=0}^n \left\| D^j \psi_{i+n-d-1} \right\|_{L^\infty}^2.$$

Using Lemma 3.1 of [3], it can be shown that there exists a positive constant  $E$  depending on  $d$  but not on  $\pi$  such that

$$\sum_{j=0}^n \left\| D^j \psi_{i+n-d-1} \right\|_{L^\infty}^2 \leq E \delta^{-2n} \quad \text{for all } 1 \leq i \leq m.$$

Hence

$$a^T R a \leq P E (2d+1) (d+1) \Delta \delta^{-2n} \|a\|_2^2,$$

and thus

$$(10) \quad \Lambda \leq P E (2d+1) (d+1) \Delta \delta^{-2n}.$$

Combining (9) and (10), we obtain the desired result with

$$C = P E (2d+1) (d+1)^2 K^{-1} D^{-2}. \quad //$$

A corollary analogous to the Corollary of [7] is clearly valid.

### Extension to higher dimensions

We consider the following class of linear, self-adjoint, boundary value problems defined on an  $M$ -dimensional hypercube  $\Omega \equiv \prod_{j=1}^M [0, 1]$  with boundary  $\partial\Omega$ :

$$(11) \quad L[u(x)] = f(x), \quad x \in \Omega, \quad f \in L^2(\Omega),$$

with homogeneous Dirichlet boundary conditions

$$(12) \quad D^\alpha u(x) = 0, \quad x \in \partial\Omega, \quad 0 \leq |\alpha| \leq n - 1, \quad n \geq 1,$$

where the linear differential operator  $L$  is defined by

$$(13) \quad L[u(x)] \equiv \sum_{0 \leq |\alpha|, |\beta| \leq n} (-1)^{|\alpha|} D^\alpha \left[ p_{\alpha\beta}(x) D^\beta u(x) \right].$$

We are using the usual multi-index notation, see [5]. Assume that all the coefficients  $p_{\alpha\beta}(x)$  are bounded measurable functions in  $\Omega$  and that

$$p_{\alpha\beta} = p_{\beta\alpha} \quad \text{for all } 0 \leq |\alpha|, |\beta| \leq n.$$

Let  $W_0^{n,2}(\Omega)$  denote the completion of the set of all  $C^\infty(\bar{\Omega})$  functions having compact support in  $\Omega$ , with respect to the Sobolev norm

$$\|w\|_{W^{n,2}} \equiv \left\{ \sum_{0 \leq |\alpha| \leq n} \int_{\Omega} [D^\alpha w(x)]^2 dx \right\}^{\frac{1}{2}}.$$

We assume that there exists a positive constant  $K$  such that for all  $w \in W_0^{n,2}(\Omega)$ ,

$$(14) \quad K \|w\|_{L^2}^2 \leq \int_{\Omega} \left\{ \sum_{0 \leq |\alpha|, |\beta| \leq n} p_{\alpha\beta}(x) D^\alpha w(x) D^\beta w(x) \right\} dx.$$

It can be shown that the problem (11)-(13) has a unique generalized solution and that the Rayleigh-Ritz method is applicable, see [6]. Let

$\{\phi_i(x)\}_{i=1}^m$  be  $m$  linearly independent coordinate functions such that

$\phi_i \in W_0^{n,2}(\Omega)$  for all  $1 \leq i \leq m$ . The Rayleigh-Ritz-Galerkin matrix

$R \equiv (r_{ik})$  for the problem (11)-(13) is given by

$$(15) \quad r_{ik} = \int_{\Omega} \left\{ \sum_{0 \leq |\alpha|, |\beta| \leq n} p_{\alpha\beta}(x) D^\alpha \phi_k(x) D^\beta \phi_i(x) \right\} dx, \quad 1 \leq i, k \leq m.$$

For each  $j$ ,  $1 \leq j \leq M$ , let  $\pi_j$  be a  $(d+1)$ -extended partition of  $[0, 1]$  in the  $j$ -th dimension:

$$\pi_j : 0 = x_0^{(j)} < x_1^{(j)} \leq x_2^{(j)} \leq \dots \leq x_{N_j}^{(j)} < x_{N_j+1}^{(j)} = 1,$$

and let  $\Delta_j$  and  $\delta_j$  be defined as in (5). Using expression (6), we

construct the normalized B-spline basis  $\left\{ \psi_k(x^{(j)}) \right\}_{k=-d+n}^{N_j-n}$  for

$Sp_0(d, \pi_j)$ ,  $j = 1, 2, \dots, M$ . Let  $\bar{\pi} \equiv \prod_{j=1}^M \pi_j$  be a  $(d+1)$ -extended product partition of  $\Omega$  and let

$$\bar{\Delta} \equiv \max_{1 \leq j \leq M} \Delta_j \quad \text{and} \quad \bar{\delta} \equiv \min_{1 \leq j \leq M} \delta_j .$$

The extended multivariate spline space  $Sp_0(d, \bar{\pi})$  is defined to be the

tensor product  $\otimes_{j=1}^M Sp_0(d, \pi_j)$ . It can be shown that  $Sp_0(d, \bar{\pi})$  is the linear span of all the normalized multivariate B-splines

$$\psi_{k_1}(x^{(1)}) \psi_{k_2}(x^{(2)}) \dots \psi_{k_M}(x^{(M)}) , \quad -d + n \leq k_j \leq N_j - n ,$$

$$j = 1, 2, \dots, M ,$$

which we rename as  $\{B_i(x)\}_{i=1}^m$ , where

$$x = (x^{(1)}, x^{(2)}, \dots, x^{(M)}) \quad \text{and} \quad m = \prod_{j=1}^M (N_j + d + 1 - 2n) .$$

Using Lemma 1, it is straightforward to prove the following:

LEMMA 2. For an arbitrary  $(d+1)$ -extended product partition  $\bar{\pi}$  of  $\Omega$ , there exists a positive constant  $\bar{D}$  depending only on  $d$  such that for all  $a \in R^m$ ,

$$(16) \quad \left\| \sum_{i=1}^m a_i B_i \right\|_{L^2} \geq \bar{D} \delta^{M/2} \|a\|_2 .$$

In applying the Rayleigh-Ritz method, let the approximating subspace  $S_m \equiv Sp_0(d, \bar{\pi})$  and let the coordinate functions  $\phi_i(x) \equiv B_i(x)$ ,  $i = 1, 2, \dots, m$ . Assuming that the maximum multiplicity of the interior knots of  $\pi_j$  is less than or equal to  $d + 1 - n$ , for all  $1 \leq j \leq M$ ,

then it can be shown that  $Sp_0(d, \bar{\pi}) \subset W_0^{n,2}(\Omega)$ , see [6].

**THEOREM 2.** *If (14) holds and  $\bar{\pi}$  is an arbitrary  $(d+1)$ -extended product partition of  $\Omega$  such that the multiplicity assumption above is valid, then there exists a positive constant  $\bar{C}$  depending only on  $d$  such that*

$$(17) \quad \kappa(R) \leq \bar{C}(\bar{\Delta}/\bar{\delta})M_{\bar{\pi}}^{-2n}.$$

**Proof.** From (15) and (14), we obtain for all  $a \in \mathbb{R}^m$ ,

$$\begin{aligned} a^T R a &= \int_{\Omega} \left\{ \sum_{0 \leq |\alpha|, |\beta| \leq n} p_{\alpha\beta}(x) D^{\alpha} \left[ \sum_{i=1}^m a_i B_i(x) \right] D^{\beta} \left[ \sum_{i=1}^m a_i B_i(x) \right] \right\} dx \\ &\geq K \left\| \sum_{i=1}^m a_i B_i \right\|_{L^2}^2, \end{aligned}$$

which, by Lemma 2, yields

$$a^T R a \geq K \bar{D}^{-2M} \|a\|_2^2,$$

and thus

$$(18) \quad \lambda \geq K \bar{D}^{-2M}.$$

Conversely, since  $p_{\alpha\beta}(x)$ ,  $0 \leq |\alpha|, |\beta| \leq n$ , are bounded in  $\Omega$ , there exists a positive constant  $Q$  such that

$$a^T R a \leq Q \sum_{0 \leq |\alpha| \leq n} \int_{\Omega} \left[ D^{\alpha} \sum_{i=1}^m a_i B_i(x) \right]^2 dx,$$

and using the minimal support properties of  $\{B_i\}_{i=1}^m$ , it can be shown that there exists a positive constant  $F$  depending only on  $d$  such that

$$\begin{aligned} a^T R a &\leq Q F \sum_{0 \leq |\alpha| \leq n} \sum_{i=1}^m a_i^2 \int_{\Omega} \left[ D^{\alpha} B_i(x) \right]^2 dx \\ &\leq Q F \sum_{i=1}^m a_i^2 (d+1)^{M_{\bar{\Delta}} M} \sum_{0 \leq |\alpha| \leq n} \left\| D^{\alpha} B_i \right\|_{L^{\infty}}^2. \end{aligned}$$

Using Lemma 3.1 of [3], it can be shown that there exists a positive constant  $\bar{E}$  depending only on  $d$  such that



$$\sum_{0 \leq |\alpha| \leq n} \left\| D^{\alpha} B_i \right\|_{L^{\infty}}^2 \leq E \delta^{-2n} \quad \text{for all } 1 \leq i \leq m .$$

Hence

$$a^T R a \leq Q \bar{F} E (d+1)^{M-M_{\Delta}} \delta^{-2n} \|a\|_2^2 ,$$

and thus

$$(19) \quad \Lambda \leq Q \bar{F} E (d+1)^{M-M_{\Delta}} \delta^{-2n} .$$

Combining (18) and (19), we obtain the desired result with

$$\bar{C} = Q \bar{F} E (d+1)^{M-K-1} \bar{D}^{-2} . \quad //$$

### References

- [1] Carl de Boor, "On uniform approximation by splines", *J. Approximation Theory* 1 (1968), 219-235.
- [2] Carl de Boor, "The quasi-interpolant as a tool in elementary polynomial spline theory", *Approximation theory*, 269-276 (Proc. Internat. Sympos., Univ. Texas, Austin, Texas, 1973. Academic Press, New York, London, 1973).
- [3] C. de Boor and G.J. Fix, "Spline approximation by quasiinterpolants", *J. Approximation Theory* 8 (1973), 19-45.
- [4] H.B. Curry and I.J. Schoenberg, "On Pólya frequency functions IV: the fundamental spline functions and their limits", *J. Analyse Math.* 17 (1966), 71-107.
- [5] С.Г. Михлин, Численная реализация вариационных методов (Izdat. "Nauka", Moscow, 1966).  
S.G. Mikhlin, *The numerical performance of variational methods* (translated by R.S. Anderssen. Wolters-Noordhoff, Groningen, 1971).
- [6] Bernard J. Omodei, "Stability of the Rayleigh-Ritz-Galerkin procedure for elliptic boundary value problems" (PhD thesis, Australian National University, Canberra, 1976). See also: Abstract, *Bull. Austral. Math. Soc.* 14 (1976), 471-472.

- [7] Martin H. Schultz, "The condition number of a class of Rayleigh-Ritz-Galerkin matrices", *Bull. Amer. Math. Soc.* 14 (1970), 840-844.

Department of Mathematics,  
University of Manchester,  
Manchester,  
England.

Present Address:  
School of Mathematical Sciences,  
Flinders University of South Australia,  
Bedford Park,  
South Australia.