

ON ERROR BOUNDS IN STRONG APPROXIMATIONS FOR EIGENVALUE PROBLEMS

R. P. KULKARNI and B. V. LIMAYE

(Received 6 March 1980)

(Revised 12 June 1980)

Abstract

Some corrections of error bounds obtained by Chatelin and Lemordant for the first three terms of the asymptotic case of a strong approximation are given. The error bounds for the approximations of order 2 in the Galerkin method are compared with the Rayleigh quotients constructed with the eigenvectors in the Sloan method. A numerical experiment is also carried out.

1. Introduction

Chatelin and Lemordant [4] have introduced the concept of a strong approximation of a bounded or unbounded closed operator T in a Banach space X and, by using the series expansions for eigenvalues and spectral projections given by Kato [7, pages 76–78], have given some error bounds for approximations of eigenvalues and spectral projections. In the present paper, we show that some of these error bounds need to be modified. As shown in [4], the analysis is relevant to most of the presently used methods for approximating differential and integral operators.

Section 2 contains the relevant notation, which agrees with that of [4]. Our major results are presented in the two theorems of Section 3. In Section 4 we give computable *a posteriori* error bounds. In this section we also compare the error bounds for the approximations of order 2 in the Galerkin method with the Rayleigh quotients constructed with the eigenvectors in the Sloan method, and

show that in the case of uniform convergence the latter give an approximation of a higher order. A numerical example at the end of the paper illustrates this comparison.

2. Preliminaries

Let X be a Banach space over C , normed by $\| \cdot \|$. Let $L(X)$ be the algebra of bounded linear operators on X and $\mathcal{C}(X)$ be the space of closed linear operators in X with domain D . Let $T \in \mathcal{C}(X)$.

The resolvent set $\rho(T)$ is the set of points z in C such that $(T - z)^{-1} \in L(X)$, the resolvent operator is $R(z) = (T - z)^{-1}$ for $z \in \rho(T)$, and the spectrum of T is $\sigma(T) = C - \rho(T)$. The resolvent operator $R(z)$ has domain X and range D .

Let T_n , $n = 1, 2, \dots$, be a sequence of operators in $\mathcal{C}(X)$. Note that, for each $z \in \rho(T)$, $(T - T_n)R(z)$ is a closed operator with domain X . By the closed graph theorem, it is bounded.

Let λ be an isolated eigenvalue of T , with finite algebraic multiplicity m and index \mathcal{C} . Let Γ be a positively oriented rectifiable curve enclosing λ , but not enclosing any other part of $\sigma(T)$. Then

$$P = -\frac{1}{2\pi i} \int_{\Gamma} R(z) dz$$

is called the spectral projection, and $M = PX$ the spectral subspace associated with λ and T . Note that the dimension of M is m .

$$S = \lim_{z \rightarrow \lambda} R(z)(1 - P)$$

is called the reduced resolvent with respect to λ .

As in Section IV of Chatelin and Lemordant [4], T_n will be called a strong approximation of T if $T_n x \rightarrow Tx$ for all $x \in D$, and $\|[(T - T_n)R(z)]^2\| \rightarrow 0$ for any z on Γ . The following special cases of strong approximation are discussed in Chatelin and Lemordant [4]: uniform approximation, collective compact approximation, neighbouring approximation, holomorphic family of type (A).

We shall assume throughout this paper that T_n is a strong approximation of T . Then it follows that, inside the curve Γ , there are exactly m eigenvalues $\lambda_{n,i}$, $i = 1, 2, \dots, m$, of T_n , counted according to algebraic multiplicities, for n large enough. (See Chatelin and Lemordant [4], Lemma 4.) We write

$$\lambda_n = \frac{1}{m} \sum_{i=1}^m \lambda_{n,i}$$

and

$$P_n = \frac{-1}{2\pi i} \int_{\Gamma} (T_n - z)^{-1} dz,$$

where P_n is the spectral projection associated with the spectrum of T_n inside Γ .

The eigenvalue λ is a pole of order l , $1 \leq l \leq m$, of $R(z)$ whose Laurent expansion can be written as follows:

$$R(z) = \sum_{k=-l}^{\infty} (z - \lambda)^k S^{(k+1)},$$

with

$$S^{(0)} = -P, \\ S^{(-k)} = -D^k, \quad k \geq 1, D = (T - \lambda)P,$$

and

$$S^{(k)} = S^k, \quad k \geq 1, S = \lim_{z \rightarrow \lambda} R(z)(1 - P).$$

Using the Neumann series of $R_n(z)$, we have (Kato [7], pages 76, 77)

$$P_n - P = \sum_{p=2}^{\infty} \sum_{\bullet} S^{(k_1)}(T - T_n)S^{(k_2)}(T - T_n)S^{(k_p)}, \tag{2.1}$$

where \bullet denotes summation over k_1, \dots, k_p subject to

$$k_1 + k_2 + \dots + k_p = p - 1$$

and

$$k_j \geq -l + 1, \quad j = 1, 2, \dots, p,$$

and

$$\lambda_n - \lambda = \frac{1}{m} \sum_{p=1}^{\infty} \frac{1}{p} \sum_{\bullet} \text{tr}[(T - T_n)S^{(k_1)} \dots (T - T_n)S^{(k_p)}]. \tag{2.2}$$

3. The main results

THEOREM 1. *Let $l = 1$, $\epsilon_n = \|(T - T_n)P\|$, $\alpha_n = \|(T - T_n)S\|$, $\beta_n = \|[(T - T_n)S]^2\|$ and $\gamma_n = \|(T - T_n)S^2(T - T_n)S\|$ so that ϵ_n and β_n tend to zero, and α_n remains bounded as $n \rightarrow \infty$. Then, for n large enough,*

(a)

$$\lambda = \lambda_n + \frac{1}{m} \text{tr}(T - T_n)P + \frac{1}{m} \text{tr}(T - T_n)S(T - T_n)P \\ + O(\max\{\beta_n \epsilon_n, \alpha_n \epsilon_n^2\}); \tag{3.1}$$

(b)

$$\|P - P_n P - S(T - T_n)P - S(T - T_n)S(T - T_n)P \\ + P(T - T_n)S^2(T - T_n)P - S^2(T - T_n)P(T - T_n)P\| \\ = O(\max\{\beta_n \epsilon_n, \alpha_n \epsilon_n^2, \gamma_n \epsilon_n, \epsilon_n^3\}). \tag{3.2}$$

PROOF. (a) Since P is a compact projection, $\epsilon_n = \|(T - T_n)P\| \rightarrow 0$ as $n \rightarrow \infty$. For the proof of $\beta_n = \|[(T - T_n)S]^2\| \rightarrow 0$ as $n \rightarrow \infty$, see Chatelin-Lemordant [4], Lemma 2, page 261. For $l = 1$, (2.2) becomes

$$\lambda_n - \lambda = \frac{1}{m} \sum_{p=1}^{\infty} \frac{1}{p} \sum_{\star} \text{tr}[(T - T_n)S^{(k_1)} \cdots (T - T_n)S^{(k_p)}],$$

where \star denotes summation over k_1, \dots, k_p subject to

$$k_1 + k_2 + \cdots + k_p = p - 1$$

and

$$k_j \geq 0, \quad j = 1, 2, \dots, p.$$

If we put $p = 1$ in the above expression, we get the term $\text{tr}(m(T - T_n)P)/m$. The remaining term,

$$\sigma = \frac{1}{m} \sum_{p=2}^{\infty} \frac{1}{p} \sum_{\star} \text{tr}[(T - T_n)S^{(k_1)} \cdots (T - T_n)S^{(k_p)}],$$

can be decomposed into the sum σ_1 over the set of all k_j 's where only one k_j is zero, the sum σ_2 over the set of k_j 's where only two k_j 's are zero, and so on. Then

$$\begin{aligned} \sigma_1 = & -\frac{1}{m} \text{tr}(T - T_n)S(T - T_n)P - \frac{1}{m} \text{tr}(T - T_n)S(T - T_n)S(T - T_n)P \\ & - \frac{1}{m} \text{tr}(T - T_n)S(T - T_n)S(T - T_n)S(T - T_n)P + \dots \end{aligned}$$

The n th term of the above expansion is

$$a_n = \begin{cases} O(\beta_n^{n/2}\epsilon_n), & \text{if } n \text{ is even} \\ O(\alpha_n\beta_n^{(n-1)/2}\epsilon_n), & \text{if } n \text{ is odd.} \end{cases}$$

It follows that

$$\sigma_1 + \frac{1}{m} \text{tr}(T - T_n)S(T - T_n)P = O(\beta_n\epsilon_n). \tag{3.3}$$

Next,

$$\begin{aligned} \sigma_2 = & \frac{1}{m} \text{tr}(T - T_n)S^2(T - T_n)P(T - T_n)P \\ & + \frac{1}{4m} [\text{tr}(T - T_n)S^2(T - T_n)S(T - T_n)P(T - T_n)P \\ & + \text{tr}(T - T_n)S^2(T - T_n)P(T - T_n)S(T - T_n)P \\ & + \text{tr}(T - T_n)S^2(T - T_n)P(T - T_n)P(T - T_n)S \\ & + \text{tr}(T - T_n)S(T - T_n)S^2(T - T_n)P(T - T_n)P \\ & + \text{tr}(T - T_n)S(T - T_n)P(T - T_n)S^2(T - T_n)P \\ & + \text{tr}(T - T_n)S(T - T_n)P(T - T_n)P(T - T_n)S^2] + \dots \end{aligned}$$

Hence it follows that $\sigma_2 = O(\alpha_n \epsilon_n^2)$ and, in general, $\sigma_k = O(\alpha_n \epsilon_n^k)$, $k \geq 2$, so that

$$\sum_{k=2}^{\infty} \sigma_k = O(\alpha_n \epsilon_n^2). \tag{3.4}$$

Thus,

$$\begin{aligned} \lambda &= \lambda_n + \frac{1}{m} \operatorname{tr}(T - T_n)P + \frac{1}{m} \operatorname{tr}(T - T_n)S(T - T_n)P \\ &\quad + O(\max\{\beta_n \epsilon_n, \alpha_n \epsilon_n^2\}). \end{aligned}$$

This proves (a).

(b) For $l = 1$, (2.1) becomes

$$P_n - P = \sum_{p=2}^{\infty} \sum_{\bullet} S^{(k_1)}(T - T_n)S^{(k_2)} \dots (T - T_n)S^{(k_p)}.$$

Since $PS = SP = 0$,

$$P_n P - P = - \sum_{p=2}^{\infty} \sum_{**} S^{(k_1)}(T - T_n)S^{(k_2)} \dots S^{(k_{p-1})}(T - T_n)P, \tag{3.5}$$

where $**$ denotes summation over k_1, \dots, k_{p-1} subject to

$$k_1 + k_2 + \dots + k_{p-1} = p - 1$$

and

$$k_j \geq 0, \quad j = 1, 2, \dots, p - 1.$$

Decompose $\sum_{p=2}^{\infty} \sum_{**} S^{(k_1)}(T - T_n)S^{(k_2)} \dots (T - T_n)P$ into the sum σ'_1 over the set of all k_j 's where only one k_j is zero, the sum σ'_2 over the set of all k_j 's where only two k_j 's are zero, and so on. Then

$$\begin{aligned} \sigma'_1 &= S(T - T_n)P + S(T - T_n)S(T - T_n)P \\ &\quad + S(T - T_n)S(T - T_n)S(T - T_n)P + \dots \end{aligned}$$

The n th term of the above expansion is

$$a'_n = \begin{cases} O(\beta^{(n-1)/2} \epsilon_n), & \text{if } n \text{ is odd} \\ O(\alpha_n \beta_n^{(n-1)/2} \epsilon_n), & \text{if } n \text{ is even.} \end{cases}$$

It follows that

$$\|\sigma'_1 - S(T - T_n)P - S(T - T_n)S(T - T_n)P\| = O(\beta_n \epsilon_n). \tag{3.6}$$

Since

$$\begin{aligned} \sigma'_2 &= -[S^2(T - T_n)P(T - T_n)P + P(T - T_n)S^2(T - T_n)P \\ &\quad + S^2(T - T_n)S(T - T_n)P(T - T_n)P + S^2(T - T_n)P(T - T_n)S(T - T_n)P \\ &\quad + S(T - T_n)S^2(T - T_n)P(T - T_n)P + S(T - T_n)P(T - T_n)S^2(T - T_n)P \\ &\quad + P(T - T_n)S(T - T_n)S^2(T - T_n)P \\ &\quad \quad \quad + P(T - T_n)S^2(T - T_n)S^2(T - T_n)S(T - T_n)P + \dots], \end{aligned}$$

we see that

$$\begin{aligned} & \| \sigma'_2 + S^2(T - T_n)P(T - T_n)P + P(T - T_n)S^2(T - T_n)P \| \\ & = O(\max\{ \alpha_n \epsilon_n^2, \beta_n \epsilon_n, \gamma_n \epsilon_n \}). \end{aligned} \tag{3.7}$$

Also,

$$\begin{aligned} \sigma'_3 & = S^3(T - T_n)P(T - T_n)P(T - T_n)P \\ & + P(T - T_n)S^3(T - T_n)P(T - T_n)P \\ & + P(T - T_n)P(T - T_n)S^3(T - T_n)P \\ & + S^3(T - T_n)S(T - T_n)P(T - T_n)P(T - T_n)P + \dots, \end{aligned}$$

so that

$$\| \sigma'_3 \| = O(\max\{ \epsilon_n^3, \alpha_n \epsilon_n^2 \}).$$

In general,

$$\| \sigma'_k \| = O(\max\{ \epsilon_n^k, \alpha_n \epsilon_n^{k-1} \}), \quad k \geq 3.$$

It follows that

$$\left\| \sum_{k=3}^{\infty} \sigma'_k \right\| = O(\max\{ \epsilon_n^3, \alpha_n \epsilon_n^2 \}). \tag{3.8}$$

From (3.5), (3.6), (3.7), and (3.8) we get

$$\begin{aligned} & \| P - P_n P - S(T - T_n)P - S(T - T_n)S(T - T_n)P \\ & + P(T - T_n)S^2(T - T_n)P + S^2(T - T_n)P(T - T_n)P \| \\ & = O(\max\{ \beta_n \epsilon_n, \alpha_n \epsilon_n^2, \gamma_n \epsilon_n \epsilon_n^3 \}). \end{aligned}$$

This proves (b).

REMARKS 1. The proof of the above theorem also shows that

$$\left. \begin{aligned} \lambda & = \lambda_n + O(\epsilon_n) \\ \| P - P_n P \| & = O(\epsilon_n) \end{aligned} \right\} \tag{3.9}$$

and

$$\left. \begin{aligned} \lambda & = \lambda_n + \frac{1}{m} \operatorname{tr}(T - T_n)P + O(\alpha_n \epsilon_n) \\ \| P - P_n P - S(T - T_n)P \| & = O(\max\{ \alpha_n \epsilon_n, \epsilon_n^2 \}). \end{aligned} \right\} \tag{3.10}$$

Let $\alpha_n \rightarrow 0$. (This is the case, for example, if $\|(T - T_n)R(z)\| \rightarrow 0$.) Then we see that error bounds in (3.10) are better than those in (3.9). Also, since $\beta_n \leq \alpha_n^2$ and $\gamma_n \leq \|S\| \alpha_n^2$, it follows that in this case the error bounds in Theorem 1 are still better than those in (3.10) above. If only $\gamma_n \rightarrow 0$ without $\alpha_n \rightarrow 0$ (which is the case, for example, if T_n is a collective compact approximation of T), then also the error bounds in Theorem 1 are better than those in (3.9) and (3.10).

2. Equation (5.2) of Chatelin and Lemordant [4] reads

$$\lambda = \lambda_n + \frac{1}{m} \operatorname{tr}(T - T_n)P + \frac{1}{m} \operatorname{tr}(T - T_n)S(T - T_n)P + O(\alpha_n \beta_n \epsilon_n).$$

In view of the proof of our equation (3.1), the above estimate is not likely to hold. A similar remark holds for the equation (5.4) of [4] in view of our equation (3.10).

3. It is stated in Chatelin and Lemordant [4], equation (5.5), that

$$\|P - P_n P + S(T - T_n)P + S(T - T_n)S(T - T_n)P - P(T - T_n)S^2(T - T_n)P\| = O(\beta_n \epsilon_n).$$

By comparing the above equation with our (3.2) of Theorem 1, it can be seen that we have added the term $-S^2(T - T_n)P(T - T_n)P$ to the expression on the left hand side of (5.5) of Chatelin and Lemordant [4]. If we do not add the above term, we get,

$$\|P - P_n P - S(T - T_n)P - S(T - T_n)S(T - T_n)P + P(T - T_n)S^2(T - T_n)P\| = O(\max\{\beta_n \epsilon_n, \gamma_n \epsilon_n, \epsilon_n^2\}).$$

The claim made in Chatelin and Lemordant [4] that the above left hand side is $O(\beta_n \epsilon_n)$ does not seem to be justified.

THEOREM 2. *Let the domain D of T be dense in X . Let $l = 1$, $\epsilon_n^* = \|(T^* - T_n^*)P^*\|$ and ϵ_n and α_n , be as in Theorem 1. Then, for n large enough,*

$$\lambda = \lambda_n + \frac{1}{m} \operatorname{tr}(T - T_n)P + \frac{1}{m} \operatorname{tr}(T - T_n)S(T - T_n)P + O(\max\{\alpha_n \epsilon_n \epsilon_n^*, \epsilon_n^2 \epsilon_n^*\}). \tag{3.11}$$

PROOF. Let x_1, x_2, \dots, x_m be a basis of $M = PX$, and $x_1^*, x_2^*, \dots, x_m^*$ be the adjoint basis. Then

$$\begin{aligned} &\operatorname{tr}(T - T_n)S(T - T_n)S(T - T_n)P \\ &= \sum_{i=1}^m (S(T - T_n)S(T - T_n)x_i, (T^* - T_n^*)x_i^*). \end{aligned}$$

Hence

$$\operatorname{tr}(T - T_n)S(T - T_n)S(T - T_n)P = O(\alpha_n \epsilon_n \epsilon_n^*).$$

Similarly,

$$\operatorname{tr}(T - T_n)S^2(T - T_n)P(T - T_n)P = O(\epsilon_n^2 \epsilon_n^*).$$

Thus we get, as in Theorem 1,

$$\lambda = \lambda_n + \frac{1}{m} \operatorname{tr}(T - T_n)P + \frac{1}{m} \operatorname{tr}(T - T_n)S(T - T_n)P + O(\max\{\alpha_n \epsilon_n \epsilon_n^*, \epsilon_n^2 \epsilon_n^*\}).$$

4. Computable error bounds

In this section we assume in addition that λ is a simple eigenvalue of T ; that is, $m = 1$. Then it can be seen that the roles of T and T_n can be interchanged, and we obtain

$$P - P_n = \sum_{p=2}^{\infty} (-1)^{p+1} \sum_{\bullet} S_n^{(k_1)}(T - T_n)S_n^{(k_2)} \cdots (T - T_n)S_n^{(k_p)}, \tag{4.1}$$

and

$$\lambda - \lambda_n = \sum_{p=1}^{\infty} \frac{(-1)^p}{p} \sum_{\bullet} \operatorname{tr}[(T - T_n)S_n^{(k_1)} \cdots (T - T_n)S_n^{(k_p)}]. \tag{4.2}$$

THEOREM 3. *Let $m = 1$. Let $\tilde{\epsilon}_n = \|(T - T_n)P_n\|$, $\tilde{\alpha}_n = \|(T - T_n)S_n\|$, $\tilde{\beta}_n = \|(T - T_n)S_n\|^2$, $\tilde{\gamma}_n = \|(T - T_n)S_n^2(T - T_n)S_n\|$. Also, let $\tilde{\epsilon}_n^* = \|(T^* - T_n^*)P_n^*\|$ in the case that D is dense in X . For n large enough,*

$$\lambda = \lambda_n + \operatorname{tr}(T - T_n)P_n + O(\tilde{\alpha}_n \tilde{\epsilon}_n) \tag{4.3}$$

$$\lambda = \lambda_n + \operatorname{tr}(T - T_n)P_n + O(\tilde{\epsilon}_n \tilde{\epsilon}_n^*) \tag{4.4}$$

$$\lambda = \lambda_n + \operatorname{tr}(T - T_n)P_n - \operatorname{tr}(T - T_n)S_n(T - T_n)P_n + O(\max\{\tilde{\beta}_n \tilde{\epsilon}_n, \tilde{\alpha}_n \tilde{\epsilon}_n^2\}) \tag{4.5}$$

$$\lambda = \lambda_n + \operatorname{tr}(T - T_n)P_n - \operatorname{tr}(T - T_n)S_n(T - T_n)P_n + O(\max\{\tilde{\alpha}_n \tilde{\epsilon}_n \tilde{\epsilon}_n^*, \tilde{\epsilon}_n^2 \tilde{\epsilon}_n^*\}) \tag{4.6}$$

$$\begin{aligned} & \|PP_n - P_n - S_n(T - T_n)P_n + S_n(T - T_n)S_n(T - T_n)P_n \\ & \quad - P_n(T - T_n)S_n^2(T - T_n)P_n - S_n^2(T - T_n)P_n(T - T_n)P_n\| \\ & = O(\max\{\tilde{\beta}_n \tilde{\epsilon}_n, \tilde{\alpha}_n \tilde{\epsilon}_n^2, \tilde{\gamma}_n \tilde{\epsilon}_n, \tilde{\epsilon}_n^3\}). \end{aligned} \tag{4.7}$$

PROOF. The proof is exactly similar to that of Theorem 1 and of Theorem 2.

REMARK. It follows from Lemma 2 and 5 of Chatelin and Lemordant [4], that $\tilde{\epsilon}_n \rightarrow 0$ and $\tilde{\beta}_n \rightarrow 0$ as $n \rightarrow \infty$. It can also be seen that $\tilde{\alpha}_n$ is bounded as $n \rightarrow \infty$ and $\tilde{\gamma}_n < \|S_n\| \tilde{\alpha}_n^2$. Further, if $T_n^* x \rightarrow T^* x$ for all $x \in \operatorname{dom}(T^*)$, then $\tilde{\epsilon}_n^* \rightarrow 0$.

Let μ_n denote the expansion of λ of order 1; that is,

$$\mu_n = \lambda_n + \text{tr}(T - T_n)P_n.$$

If ϕ_n is a normalized eigenvector of T_n associated with λ_n and ϕ_n^* is an eigenvector of T_n^* associated with $\bar{\lambda}_n$ satisfying $(\phi_n, \phi_n^*) = 1$, then $\mu_n = \text{tr} TP_n = (T\phi_n, \phi_n^*)$. From (4.3) and (4.4) we have

$$\left. \begin{aligned} \lambda &= \mu_n + O(\tilde{\alpha}_n \tilde{\epsilon}_n), \\ \lambda &= \mu_n + O(\tilde{\epsilon}_n^* \tilde{\epsilon}_n). \end{aligned} \right\} \tag{4.8}$$

If ν_n denotes the expansion of λ of order 2, that is,

$$\nu_n = \lambda_n + \text{tr}(T - T_n)P_n - \text{tr}(T - T_n)S_n(T - T_n)P_n,$$

then it follows from (4.5) and (4.6) that

$$\left. \begin{aligned} \lambda &= \nu_n + O(\max\{\tilde{\beta}_n \tilde{\epsilon}_n, \tilde{\alpha}_n \tilde{\epsilon}_n^2\}) \\ \lambda &= \nu_n + O(\max\{\tilde{\alpha}_n \tilde{\epsilon}_n \tilde{\epsilon}_n^*, \tilde{\epsilon}_n^2 \tilde{\epsilon}_n^*\}). \end{aligned} \right\} \tag{4.9}$$

If $\tilde{\epsilon}_n \rightarrow 0$, but $\tilde{\epsilon}_n^* \not\rightarrow 0$, from (4.3) and (4.5) we see that ν_n improves on μ_n .

If both $\tilde{\epsilon}_n$ and $\tilde{\epsilon}_n^*$ tend to zero, it is clear from (4.4) and (4.6) that ν_n and μ_n are approximations of the order of $O(\tilde{\epsilon}_n \tilde{\epsilon}_n^*)$. Thus ν_n does not always improve on μ_n .

If $\tilde{\alpha}_n \rightarrow 0$, irrespective of the behaviour of $\tilde{\epsilon}_n^*$, ν_n improves on μ_n .

Chatelin has considered the following special case in [5]: $T: X \rightarrow X$ is a bounded linear operator, and $\pi_n, n = 1, 2, \dots$, is a sequence of projections of X on a family of finite dimensional subspaces $X_n = \pi_n X$. It is assumed that $\pi_n x \rightarrow x, x \in X$. Then T is approximated by $T_n^G = \pi_n T \pi_n$, which is Galerkin's method. For compact T then T_n^G is a collective compact approximation of T . If T is not compact, it is assumed that T_n^G is a strong approximation of T . Let $\tilde{\epsilon}_n^G = \|(T - T_n^G)P_n^G\|$ and let $\tilde{\beta}_n^G$ and $\tilde{\alpha}_n^G$ be defined similarly. It is claimed in Chatelin [5], page 1120, that

$$|\lambda - \mu_n^G| = O(\tilde{\epsilon}_n^G \tilde{\epsilon}_n^{*G}),$$

and

$$|\lambda - \nu_n^G| = \left\{ \begin{aligned} &O(\tilde{\epsilon}_n^G \tilde{\epsilon}_n^{*G}) \quad \text{in general,} \\ &O(\|(1 - \pi_n)T\| \tilde{\epsilon}_n^G \tilde{\epsilon}_n^{*G}) \quad \text{if } T \text{ is compact.} \end{aligned} \right\}$$

(In the notations used there, $\nu_n^G = \lambda_n^{1S} = \lambda_n^{2G}, \mu_n^G = \lambda_n^G$.)

As we can see from (4.9) above, when T is bounded but not compact, the error bound for $|\lambda - \nu_n^G|$ should be modified to

$$\begin{aligned} |\lambda - \nu_n^G| &= O(\max\{\tilde{\beta}_n^G \tilde{\epsilon}_n^G, \tilde{\alpha}_n^G (\tilde{\epsilon}_n^G)^2\}) \\ &= O(\max\{\tilde{\alpha}_n^G \tilde{\epsilon}_n^G \tilde{\epsilon}_n^{*G}, (\tilde{\epsilon}_n^G)^2 \tilde{\epsilon}_n^{*G}\}) \\ &= O(\tilde{\epsilon}_n^G \tilde{\epsilon}_n^{*G}). \end{aligned}$$

It is said in Chatelin [5] that ν_n^G always improves on $\mu_n^G = \lambda_n^G$ as a consequence of the stepwise convergence. But stepwise convergence is proved there only for the collectively compact convergence and not for the strong convergence.

We now consider an approximation of the order of $\tilde{\epsilon}_n^2$. Let T be a (densely defined) self-adjoint operator on a Hilbert space H . If $\rho_n = (T\phi_n, \phi_n)$ denote the *Rayleigh quotients*, it can be seen easily that $\rho_n \rightarrow \lambda$ as $n \rightarrow \infty$. In fact, if we let d_n to be the distance of ρ_n from $\sigma(T) - \{\lambda\}$, then

$$|\lambda - (T\phi_n, \phi_n)| \leq \frac{1}{d_n} [\|(T - T_n)\phi_n\|^2 - |((T - T_n)\phi_n, \phi_n)|^2].$$

This can be seen as follows:

$$\begin{aligned} & \|(T - T_n)\phi_n\|^2 - |((T - T_n)\phi_n, \phi_n)|^2 \\ &= ((T - T_n)\phi_n, (T - T_n)\phi_n) - |(T\phi_n, \phi_n) - \lambda_n|^2 \\ &= \|T\phi_n\|^2 - (\lambda_n\phi_n, T\phi_n) - (T\phi_n, \lambda_n\phi_n) + |\lambda_n|^2 - |\rho_n - \lambda_n|^2 \\ &= \|T\phi_n\|^2 - 2\rho_n \operatorname{Re} \lambda_n + |\lambda_n|^2 - \rho_n^2 - |\lambda_n|^2 + 2\rho_n \operatorname{Re} \lambda_n \\ &= \|T\phi_n\|^2 - \rho_n^2. \end{aligned}$$

Hence, by using the Kato-Temple inequality [6],

$$\begin{aligned} |\lambda - (T\phi_n, \phi_n)| &\leq \frac{1}{d_n} [\|T\phi_n\|^2 - \rho_n^2] \\ &= \frac{1}{d_n} [\|(T - T_n)\phi_n\|^2 - |((T - T_n)\phi_n, \phi_n)|^2]. \end{aligned}$$

Since the d_n 's are bounded away from zero, we see that

$$\lambda = \rho_n + O(\tilde{\epsilon}_n^2). \tag{4.10}$$

Let us now consider the following two particular kinds of approximations: T is a bounded self-adjoint operator on a Hilbert space H and $\pi_n: H \rightarrow H_n$ is a projection, where H_n is a closed subspace of H . Let $T_n^G = \pi_n T \pi_n$ be the approximating operator in *Galerkin's method*. In *Sloan's method*, T is approximated by $T_n^S = T \pi_n$ (Sloan [8]). Assume that T_n^G and T_n^S are strong approximations of T (which is the case, for example, if T is compact, and $\pi_n x \rightarrow x$, $x \in H$). If ϕ_n^S denotes a normalized eigenvector of T_n^S associated with λ_n^S and $\rho_n^S = (T\phi_n^S, \phi_n^S)$, then (4.10) shows that

$$\lambda = \rho_n^S + O(\|(T - T_n^S)P_n^S\|^2).$$

But

$$\begin{aligned} \|(T - T_n^S)P_n^S\| &= \|T(1 - \pi_n)P_n^S\| \\ &= \|T(1 - \pi_n)^2 P_n^S\| \\ &\leq \|T(1 - \pi_n)\| \|(1 - \pi_n)P_n^S\|. \end{aligned}$$

Also,

$$\|(1 - \pi_n)P_n^S\| = O(\|(T - T_n^G)P_n^G\|),$$

as proved on page 1118 of Chatelin [5].

If we let $\tilde{\epsilon}_n^G = \|(T - T_n^G)P_n^G\|$, and similarly for $\tilde{\beta}_n^G$ and $\tilde{\alpha}_n^G$, then

$$\begin{aligned} \lambda &= \rho_n^S + O(\|T - T\pi_n\|^2(\tilde{\epsilon}_n^G)^2) \\ &= \rho_n^S + o(\tilde{\epsilon}_n^G) \end{aligned} \tag{4.11}$$

This should be compared with

$$\begin{aligned} \lambda &= \nu_n^G + O(\max\{\tilde{\beta}_n^G\tilde{\epsilon}_n^G, \tilde{\alpha}_n^G(\tilde{\epsilon}_n^G)^2\}) \\ &= \nu_n^G + o(\tilde{\epsilon}_n^G), \end{aligned} \tag{4.12}$$

as obtained from (4.9). Thus, ρ_n^S and ν_n^G are both approximations of λ of the order of $o(\tilde{\epsilon}_n^G)$.

Further, since $P_n^G H \subset H_n$, we see that

$$\tilde{\epsilon}_n^G = \|(T - \pi_n T \pi_n)P_n^G\| = O(\|T - \pi_n T\|).$$

Hence, by (4.11) and (4.12),

$$\lambda = \rho_n^S + O(\|T - T\pi_n\|^2\|T - \pi_n T\|^2), \tag{4.13}$$

and

$$\lambda = \nu_n^G + O(\max\{\|T - \pi_n T \pi_n\|^2\|T - \pi_n T\|, \|T - \pi_n T \pi_n\| \|T - \pi_n T\|^2\}). \tag{4.14}$$

If each π_n is an orthogonal projection, then $\tilde{\epsilon}_n^G = \hat{\epsilon}_n^{*G}$ so that, by (4.13) and (4.9),

$$\lambda = \rho_n^S + O(\|T - \pi_n T\|^4) \tag{4.15}$$

and

$$\lambda = \nu_n^G + O(\max\{\|T - \pi_n T \pi_n\| \|T - \pi_n T\|^2, \|T - \pi_n T\|^3\}). \tag{4.16}$$

In the case that the Galerkin approximations and the Sloan approximations are uniform, that is, if $\|\pi_n T \pi_n - T\| \rightarrow 0$ and $\|T \pi_n - T\| \rightarrow 0$ (which is the case, for example, if T is compact, $\pi_n x \rightarrow x$ and $\pi_n^* x \rightarrow x$), then the results (4.13) and (4.14) give comparative estimates for the two approximations ρ_n^S and ν_n^G of λ . In particular, if each π_n is orthogonal, then (4.15) and (4.16) show that ρ_n^S is a better approximation than ν_n^G . This is borne out by the following example.

5. A numerical example

Let $X = L^2([0, 1])$, and T be the operator defined by

$$Tx(s) = \int_0^1 k(s, t)x(t) dt,$$

where

$$k(s, t) = \begin{cases} (1 - t)s, & \text{if } t \geq s \\ (1 - s)t, & \text{if } t < s. \end{cases}$$

Then T is a self-adjoint compact operator on $L^2([0, 1])$. The eigenvalues of T are given by $\lambda^k = 1/(k^2\pi^2)$, $k = 1, 2, \dots$. Consider the following orthonormal basis in $L^2([0, 1])$:

$$e_1(s) = 1, \quad e_k(s) = 2^{1/2} \cos(k - 1)\pi s, \quad k = 2, 3, \dots$$

Let $X_n = \text{span}\{e_1, e_2, \dots, e_n\}$ and $\pi_n: X \rightarrow X_n$ denote the orthogonal projection of X onto X_n . The matrix corresponding to the Galerkin approximation $\pi_n T \pi_n$ is

$$A_n = (Te_j, e_i), \quad i, j = 1, \dots, n.$$

The matrix A_n is real and symmetric with

$$(A_n)_{i,j} = \begin{cases} \frac{1}{12} & \text{if } i = j = 1, \\ \frac{-\sqrt{2}}{(j - 1)^2 \pi^2} & \text{if } i = 1, j > 2, j \text{ odd}, \\ 0 & \text{if } i = 1, j > 2, j \text{ even}, \\ \frac{1}{(j - 1)^2 \pi^2} & \text{if } i = j, i \text{ odd}, \\ \frac{1}{(j - 1)^2 \pi^2} - \frac{8}{(j - 1)^4 \pi^4} & \text{if } i = j, i \text{ even}, \\ 0 & \text{if } i \neq j, \text{ either } i \text{ or } j \text{ is odd}, \\ -\frac{8}{(j - 1)^2 (i - 1)^2 \pi^4} & \text{if } i \neq j, \text{ both } i \text{ and } j \text{ are even.} \end{cases}$$

If $(\phi_n)_i$, $i = 1, 2, \dots, n$, is an eigenvector for A_n corresponding to an eigenvalue λ_n , then $\phi_n = \sum_{i=1}^n (\phi_n)_i e_i$ is an eigenvector of $\pi_n T \pi_n$ associated with the eigenvalue λ_n , and $T\phi_n = \sum_{i=1}^n (\phi_n)_i Te_i$ is an eigenvector of $T\pi_n$ associated with λ_n . The Rayleigh quotient $\rho_n^S = (T\phi_n^S, \phi_n^S) / (\phi_n^S, \phi_n^S)$ where $\|\phi_n^S\| = 1$ is then $(T^2\phi_n, T\phi_n) / (T\phi_n, T\phi_n)$.

Also, it follows, since $\phi_n^* = \phi_n$, that $\nu_n^G = (T^2\phi_n, \phi_n) / \lambda_n$ with $\sum_{i=1}^n (\phi_n)_i^2 = 1$. (See page 269 of Chatelin and Lemordant [4].)

To compute ρ_n^S and ν_n^G , we note that

$$T\phi_n = \sum_{i=1}^n (\phi_n)_i Te_i,$$

where

$$Te_i(s) = \begin{cases} \frac{s(1-s)}{2} & \text{if } i = 1, \\ \left(-\frac{\sqrt{2}}{(i-1)^2\pi^2}\right)e_1 + \left(\frac{1}{(i-1)^2\pi^2}\right)e_i, & \text{if } i \geq 2, i \text{ odd,} \\ \left(-\frac{\sqrt{2}}{(i-1)^2\pi^2}\right)e_1 + \left(\frac{1}{(i-1)^2\pi^2}\right)e_i + \left(\frac{2\sqrt{2}}{(i-1)^2\pi^2}\right)s, & \text{if } i \geq 2, i \text{ even.} \end{cases}$$

Hence

$$T^2\phi_n = \sum_{i=1}^n (\phi_n)_i T^2e_i,$$

where

$$T^2e_i(s) = \begin{cases} \frac{s(1-s^2)^2}{24}, & \text{if } i = 1, \\ \left(-\frac{\sqrt{2}}{(i-1)^2\pi^2}\right)Te_1 + \left(\frac{1}{(i-1)^2\pi^2}\right)Te_i, & \text{if } i \geq 2, i \text{ odd,} \\ \left(-\frac{\sqrt{2}}{(i-1)^2\pi^2}\right)Te_1 + \left(\frac{1}{(i-1)^2\pi^2}\right)Te_i \\ + \left(\frac{2\sqrt{2}}{(i-1)^2\pi^2}\right)\left(\frac{s(1-s^2)}{6}\right), & \text{if } i \geq 2, i \text{ even.} \end{cases}$$

We compute, in double precision, the first four eigenvalues $\lambda_n^1, \dots, \lambda_n^4$ of the matrix A_n and the corresponding eigenvectors for $n = 4, 8, 12, 16, 20$ and 30 , and obtain the corresponding values of ρ_n^S and ν_n^G . In Table 1 the values of $\lambda^1 - \lambda_n^1, \lambda^1 - \nu_n^1, \lambda^4 - \lambda_n^4$ and $\lambda^4 - \nu_n^4$ are taken from Chatelin and Lemordant [4], while other values are computed by us using the above method. As noted earlier, ρ_n^S turns out to be a better approximation of λ than ν_n^G since the Galerkin and the Sloan approximations converge uniformly in this case and the projections π_n are orthogonal.

TABLE 1
Calculations for the numerical example

n	4	8	12	16	20	30
$\lambda^1 - \lambda_n^1$	9×10^{-4}	7.8×10^{-5}	2×10^{-5}	8×10^{-6}	3.9×10^{-6}	1.1×10^{-6}
$\lambda^1 - \nu_n^1$	4.2×10^{-5}	8.9×10^{-7}	1.2×10^{-7}	4.8×10^{-8}	2.2×10^{-8}	1.7×10^{-8}
$\lambda^1 - \rho_n^{S1}$	4×10^{-6}	2.6×10^{-8}	1.5×10^{-9}	2.2×10^{-10}	5.1×10^{-11}	3.9×10^{-2}
$\lambda^2 - \lambda_n^2$	4.5×10^{-4}	5.4×10^{-5}	1.6×10^{-5}	6.7×10^{-6}	3.4×10^{-6}	1×10^{-6}
$\lambda^2 - \nu_n^2$	5.4×10^{-5}	1.9×10^{-6}	2.4×10^{-7}	4.6×10^{-8}	3.2×10^{-8}	1.5×10^{-8}
$\lambda^2 - \rho_n^{S2}$	1.2×10^{-5}	1.6×10^{-7}	1.2×10^{-8}	2×10^{-9}	4.8×10^{-10}	3.9×10^{-11}
$\lambda^3 - \lambda_n^3$	3×10^{-3}	8.6×10^{-5}	2.1×10^{-5}	8.2×10^{-6}	4×10^{-6}	1.1×10^{-6}
$\lambda^3 - \nu_n^3$	6.4×10^{-4}	8.7×10^{-6}	8.9×10^{-7}	2×10^{-7}	6.2×10^{-8}	1.2×10^{-8}
$\lambda^3 - \rho_n^{S3}$	-2.2×10^{-5}	1.2×10^{-6}	5.5×10^{-8}	6.9×10^{-9}	1.4×10^{-9}	8.7×10^{-11}
$\lambda^4 - \lambda_n^4$	1.7×10^{-3}	6.1×10^{-5}	1.7×10^{-5}	6.9×10^{-6}	3.5×10^{-6}	1×10^{-6}
$\lambda^4 - \nu_n^4$	4.3×10^{-4}	8.7×10^{-6}	1.1×10^{-6}	2.7×10^{-7}	9.8×10^{-8}	2.6×10^{-8}
$\lambda^4 - \rho_n^{S4}$	4×10^{-5}	1.8×10^{-6}	1.1×10^{-7}	1.6×10^{-8}	3.7×10^{-9}	2.7×10^{-10}

Acknowledgement

This work was supported in part by a grant from the Council of Scientific and Industrial Research, New Delhi, India.

References

- [1] P. M. Anselone, *Collectively compact operator approximation theory* (Prentice-Hall, New Jersey, 1971).
- [2] F. Chatelin and J. Lemordant, "Error bounds in the approximation of eigenvalues of differential and integral operators", *Report No. STAN-CS-75-479, Stanford University* (1975).
- [3] F. Chatelin, *Linear spectral approximation in Banach spaces* (Academic Press, New York), to appear.
- [4] F. Chatelin and J. Lemordant, "Error bounds in the approximation of eigenvalues of differential and integral operators", *J. Math. Anal. and Appl.* 62 (1978), 257–271.
- [5] F. Chatelin, "Numerical computation of the eigenelements of linear integral operators by iterations", *SIAM J. Numer. Anal.* 15 (1978), 1112–1124.
- [6] T. Kato, "On the upper and lower bounds for eigenvalues", *J. Phys. Soc. Japan* 4 (1949), 334–339.
- [7] T. Kato, *Perturbation theory for linear operators* (Springer-Verlag, Berlin, 1966).
- [8] I. H. Sloan, "Iterated Galerkin method for eigenvalue problems", *SIAM J. Numer. Anal.* 13 (1976), 753–760.

Department of Mathematics
Indian Institute of Technology
Powai
Bombay 400076
India