



# Guaranteed error estimate for solutions to two-point boundary value problems

Akitoshi Takayasu<sup>†</sup>, Shin'ichi Oishi<sup>‡</sup> and Takayuki Kubo<sup>§</sup>

<sup>†</sup>Graduate School of Fundamental Science and Engineering, Waseda University  
Okubo 3-4-1, Shinjuku-ku, Tokyo 169-8555 Japan  
<sup>‡</sup>Faculty of Science and Engineering, Waseda University  
Okubo 3-4-1, Shinjuku-ku, Tokyo 169-8555 Japan  
<sup>§</sup>Institute of Mathematics, University of Tsukuba  
Tenodai 1-1-1, Tsukuba, Ibaraki 305-8571 Japan  
Email: takitoshi@suou.waseda.jp, oishi@waseda.jp, tkubo@math.tsukuba.ac.jp

**Abstract**—We consider the ‘guaranteed’ error estimate for solutions to two-point boundary value problems. ‘Guaranteed error’ requires every error in solving the system. It is important and difficult to analyze the error which is caused by the truncation or the discretization of problems. We overcome these difficulties by using Newton-Kantorovich theorem and the solution operator to linearized problem. The original problem is transformed into a nonlinear operator equation. The guaranteed error is bounded by Newton-Kantorovich theorem through verifying some constants concerning the solution operator. Finally, numerical result is presented.

## 1. Introduction

This article is concerned with two-point boundary value problems of the form:

$$\begin{cases} -u'' = ru^N + f & 0 < x < 1, \\ u(0) = u(1) = 0, \end{cases} \quad (1)$$

where  $N \geq 2$  is a natural number,  $r \in L^\infty([0, 1])$  and  $f \in L^2([0, 1])$ . We propose a numerical verification method to prove the existence of solutions to problem (1). The goal of our verification method is to get the guaranteed error estimate. It is bounded by the following form:

$$\|u - \hat{u}\|_X \leq Const.$$

where  $u$  is the exact solution,  $\hat{u}$  is an approximate solution,  $X$  is a suitable functional space and  $Const.$  is computable. The guaranteed error estimate is rigorous i.e. it includes all computational error such as the discretization error and the rounding error when solving the problems. Namely, we can solve the two-point boundary value problem with mathematically rigorous by our verification method.

The main point of the verification method is to transform the problem (1) into a nonlinear operator equation with a solution operator to the following linearized problem:

$$\begin{cases} -u'' = f & 0 < x < 1, \\ u(0) = u(1) = 0. \end{cases} \quad (2)$$

In [4, 6], there are several methods which have been developed to bound the guaranteed error when  $N = 2$ . The solution operator is denoted as an integral operator by Green function in [4]. (1) is transformed into the integral equation. The approximate solution operator with discretizing the integral equation is used. We modified this operator argument in finite element method in [6]. An approximate solution operator is defined as a matrix form. The solution operator is estimated by the minimal eigenvalue of the corresponding problem to (2) and some operator formulation. We applied Banach’s fixed point theorem to fixed point formulation by the solution operator. Then, the guaranteed error estimate is given. However, it is difficult to formulate the fixed-point formulation if  $N \geq 3$ .

In this article, Newton-Kantorovich theorem is applied to a nonlinear operator equation corresponding to (1). We check some conditions of the theorem concerning the nonlinear operator. There are three constants: the estimation of the inverse operator, the residual of the operator equation and the Lipschitz constant of the Fréchet derivatives. The inverse operator estimation is given by Theorem2. If three assumptions of the theorem are obtained, we can estimate the norm of the inverse operator. The norm estimation gives us the residual of the operator equation. We have the upper bound of the residual. Moreover, the Lipschitz constant needs the condition that the parameter is in the open ball centered at the approximate solution with some radius. By using this condition, we have the Lipschitz constant.

We introduce our verification procedure in Section 2. First the solution operator is defined. We explain the discretization of the problem. The approximate solution operator is also defined. After that we show Newton-Kantorovich theorem. In order to obtain the assumptions of Newton-Kantorovich theorem, we need three constants. We compute these constants by using the argument of the solution operator. Finally a numerical result is presented in Section 3.

## 2. Numerical Verification Method

In this section, we shall propose the following verification procedure. In order to get the guaranteed error estimate

between the exact solution and an approximate solution by finite element method, Newton-Kantorovich theorem is applied to nonlinear operator equation. We first introduce the solution operator of finite element method.

## 2.1. Solution operator and its approximate operator

A solution operator  $\mathcal{K}$  leads the solution  $u = \mathcal{K}f$  of linearized equation (2). Here, in this section, we set  $X = L^2([0, 1])$  as the functional space. The solution operator is linear bounded operator:  $\mathcal{K} \in \mathcal{L}(X, X)$ . We note that  $\mathcal{K} : X \rightarrow H_0^1 \subset X$  is the compact operator by Sobolev embedding theorem.

After that we consider an approximate operator of  $\mathcal{K}$ . Let

$$\Delta : 0 < x_1 < \dots < x_n < 1$$

be a partition of the interval  $[0, 1]$ . Discretization is  $n + 1$  divide equally on  $[0, 1]$ .  $h$  is the width of an interval  $[x_{i-1}, x_i]_{i=2, \dots, n}$ .  $\mathcal{P}_n : X \rightarrow X_n$  is defined as a discrete projection. Let  $S_h$  be the class of basis in finite element method. We denote that a discrete functional space of  $X$  is  $X_n \subset H_0^1$ , which is defined as the set of polynomials that satisfies the boundary conditions. By these bases, a discrete functional space is composed

$$X_n = \text{span}\{\phi_{h_1}, \phi_{h_2}, \dots, \phi_{h_n}\}, \quad \phi_{h_i} \in S_h.$$

An element of  $X_n$  is described

$$u_h \in X_n, \quad u_h = \sum_{j=1}^n u_j \phi_{h_j}$$

where  $u_j = u(x_j)$ . Moreover, we choose a norm of  $X_n$

$$\|u_h\|_{X_n} = \|u_h\|_2 = \sqrt{\sum_{j=1}^n u_j^2}.$$

According to the standard argument of finite element method [2], the linearized problem (2) is transformed into the following weak formula:

$$(u', \phi'_h) = (f, \phi_h), \quad \phi_h \in X_n.$$

Here setting

$$u_h = \sum_{j=1}^n u_j \phi_{h_j}, \quad f_h = \sum_{j=1}^n f_j \phi_{h_j},$$

$$D_n = (\phi'_{h_j}, \phi'_{h_i})_{i,j=1}^n, \quad A_n = (\phi_{h_j}, \phi_{h_i})_{i,j=1}^n.$$

Discretizing weak formula, we obtain the finite linear system:

$$D_n U_h = A_n F_h$$

where  $U_h = (u_1, \dots, u_n)^T$ ,  $F_h = (f_1, \dots, f_n)^T$  are coefficients of discrete function  $u_h$  and  $f_h$ . If  $D_n$  has the inverse matrix, a finite element solution is written as follows:

$$U_h = D_n^{-1} A_n F_h.$$

On the other hand, an expression of the solution operator says that  $u = \mathcal{K}f$ . Discretizing  $f$  and  $\mathcal{K}$ , the approximate solution is obtained

$$u_h = \mathcal{P}_n \mathcal{K} f_h.$$

By identifying  $U_h$  with  $u_h$  and  $F_h$  with  $f_h$ , the approximate solution operator  $\mathcal{P}_n \mathcal{K} : X_n \rightarrow X_n$  is defined as the matrix form:

$$\mathcal{P}_n \mathcal{K} = D_n^{-1} A_n. \quad (3)$$

## 2.2. Newton-Kantorovich theorem

By using the solution operator, the problem is transformed into operator equation:

$$(1) \iff u = \mathcal{K}(ru^N + f).$$

Nonlinear operator equation is defined

$$F(u) = u - \mathcal{K}ru^N - \mathcal{K}f = 0. \quad (4)$$

Here, we assume that we can find a finite element solution, which is a good approximation of the solution. Then, the next step of our method is to check conditions of Newton-Kantorovich theorem:

**Theorem 1 (Newton-Kantorovich Theorem)** *Let  $F$  be a nonlinear operator defined by (4). We assume that the Fréchet derivative  $F'(u)$  is nonsingular and satisfies the inequality:*

$$\|F'(\hat{u})^{-1} F(\hat{u})\|_X \leq \alpha,$$

for a certain positive constant  $\alpha$ , where  $\hat{u}$  is an approximate solution to (4). Furthermore, we assume that  $F$  satisfies

$$\|F'(\hat{u})^{-1} (F'(v) - F'(w))\|_X \leq \omega \|v - w\|_X$$

with a certain positive constant  $\omega$ , for  $\forall v, w \in B(\hat{u}, \delta) \subset X$ , which is an open ball centered at  $\hat{u}$  with radius  $\delta$ . If

$$\alpha\omega \leq \frac{1}{2} \quad (5)$$

and

$$\rho = \frac{1 - \sqrt{1 - 2\alpha\omega}}{\omega},$$

then there exists the unique exact solution  $u$  to (4) in  $\bar{B}(\hat{u}, \rho)$ . Therefore the guaranteed error estimate is given by

$$\|u - \hat{u}\|_X \leq \rho. \quad (6)$$

In this way, we can show the existence and the uniqueness of the exact solution. Additionally, the guaranteed error is bounded by Newton-Kantorovich theorem. In order to verify the theorem, we need to compute three constants  $C_1, C_2$  and  $C_3$ . These satisfy

$$\|F'(\hat{u})^{-1}\|_{\mathcal{L}(X, X)} = \|(I - N\mathcal{K}r\hat{u}^{N-1})^{-1}\|_{\mathcal{L}(X, X)} \leq C_1,$$

$$\|F(\hat{u})\|_X = \|\hat{u} - \mathcal{K}r\hat{u}^N - \mathcal{K}f\|_X \leq C_2$$

and

$$\begin{aligned} & \|F'(v) - F'(w)\|_X \\ &= \|N\mathcal{K}r(v^{N-2} + v^{N-3}w + \dots + vw^{N-3} + w^{N-2})(v-w)\|_X \\ &\leq C_3\|v-w\|_X. \end{aligned}$$

Hence, we set  $\alpha = C_1 * C_2$ ,  $\omega = C_1 * C_3$ . If three constants are estimated and the condition (5) is obtained, then we get the guaranteed error estimate in the form (6).

### 2.3. Some constants

In this part, we shall concern with some constants regarding Newton-Kantorovich theorem. We explain some techniques to compute these constants. First of all,  $C_1$  is estimated according to the following theorem with respect to the inverse operator.

**Theorem 2 (Oishi[4] 7.2)** *Let  $\mathcal{K} : X \rightarrow X$  be the compact operator and  $\mathcal{P}_n : X \rightarrow X_n$  be the projection operator where  $X_n \subset X$  is the discrete functional space on  $X$ . We assume that  $\mathcal{P}_n\mathcal{K}$  is estimated as*

$$\|\mathcal{P}_n\mathcal{K}\|_{\mathcal{L}(X,X_n)} \leq K,$$

a difference between  $\mathcal{K}$  and  $\mathcal{P}_n\mathcal{K}$  is estimated as

$$\|\mathcal{K} - \mathcal{P}_n\mathcal{K}\|_{\mathcal{L}(X,X)} \leq L$$

and finite dimensional operator  $(I - \mathcal{P}_n\mathcal{K}) : X_n \rightarrow X_n$  has the inverse operator, which satisfies

$$\|(I - \mathcal{P}_n\mathcal{K})^{-1}\|_{\mathcal{L}(X_n,X_n)} \leq M.$$

If three assumptions are obtained and  $(1 + MK)L < 1$ , then the operator  $(I - \mathcal{K})$  has the inverse operator and that is estimated as

$$\|(I - \mathcal{K})^{-1}\|_{\mathcal{L}(X,X)} \leq \frac{1 + MK}{1 - (1 + MK)L}.$$

In order to compute the constant  $C_1$ , three constants  $K_1, L_1$  and  $M_1$  is needed in assumptions of Theorem2. Since  $\mathcal{K}$  is compact operator,  $r \in L^\infty([0, 1])$  and  $\hat{u} \in X_n$ , we see the operator  $N\mathcal{K}r\hat{u}^{N-1}$  is compact. Then, we have

$$\|\mathcal{P}_nN\mathcal{K}r\hat{u}^{N-1}\|_{\mathcal{L}(X,X_n)} \leq N\|\mathcal{P}_n\mathcal{K}\|_{\mathcal{L}(X,X_n)}\|r\|_\infty\|\hat{u}^{N-1}\|_\infty \leq K_1,$$

and

$$\begin{aligned} & \|N\mathcal{K}r\hat{u}^{N-1} - \mathcal{P}_nN\mathcal{K}r\hat{u}^{N-1}\|_{\mathcal{L}(X,X)} \\ &\leq N\|\mathcal{K} - \mathcal{P}_n\mathcal{K}\|_{\mathcal{L}(X,X)}\|r\|_\infty\|\hat{u}^{N-1}\|_\infty \\ &\leq L_1. \end{aligned}$$

To be more precise, the estimation of  $\|\mathcal{P}_n\mathcal{K}\|_{\mathcal{L}(X,X_n)} \leq K$  and  $\|\mathcal{K} - \mathcal{P}_n\mathcal{K}\|_{\mathcal{L}(X,X)} \leq L$  plays important role in the verification method. These are given as follows:

$\|\mathcal{P}_n\mathcal{K}\|_{\mathcal{L}(X,X_n)} \leq K$  is obtained by the minimal eigenvalue. The eigenvalue problem:

$$\begin{cases} -u'' = \lambda u & 0 < x < 1, \\ u(0) = u(1) = 0, \end{cases}$$

has eigenvalues  $\lambda = n^2\pi^2$ , then the minimal eigenvalue is  $\lambda_{min} = \pi^2$ . We have the estimation of  $\mathcal{K}$

$$\begin{aligned} \|\mathcal{K}\|_{\mathcal{L}(X,X)} &= \sup_{f \in X} \frac{\|\mathcal{K}f\|_X}{\|f\|_X} \\ &\leq \sup_{f \in X} \frac{\lambda_{min}^{-1}\|u''\|_X}{\|f\|_X} = \frac{1}{\pi^2}. \end{aligned}$$

According to the argument of the discrete projection  $X$  to  $X_n$ ,  $\|\mathcal{P}_n\|_{\mathcal{L}(X,X_n)} \leq 1$ . Thus, we get the constant  $K$

$$\begin{aligned} \|\mathcal{P}_n\mathcal{K}\|_{\mathcal{L}(X,X_n)} &\leq \|\mathcal{P}_n\|_{\mathcal{L}(X,X_n)}\|\mathcal{K}\|_{\mathcal{L}(X,X)} \\ &\leq \frac{1}{\pi^2} = K. \end{aligned}$$

$\|\mathcal{K} - \mathcal{P}_n\mathcal{K}\|_{\mathcal{L}(X,X)} \leq L$  is given by the error estimate of FEM. By using the error estimation of FEM (Aubin-Nitshe's trick):

$$\|u - \mathcal{P}_nu\|_X \leq \frac{h^2}{\pi^2}\|u''\|_X,$$

we have

$$\begin{aligned} \|\mathcal{K} - \mathcal{P}_n\mathcal{K}\|_{\mathcal{L}(X,X)} &\leq \sup_{f \in X} \frac{\|\mathcal{K}f - \mathcal{P}_n\mathcal{K}f\|_X}{\|f\|_X} \\ &= \sup_{f \in X} \frac{\|u - \mathcal{P}_nu\|_X}{\|f\|_X} \\ &\leq \frac{h^2}{\pi^2} = L. \end{aligned}$$

Furthermore, by using the approximate solution operator  $\mathcal{P}_n\mathcal{K}$ , we get the estimation of  $M_1$ . Let  $B_n$  be the matrix:

$$B_n = (Nr(x)\hat{u}^{N-1}(x)\phi_{h_j}, \phi_{h_i})_{i,j=1}^n.$$

Assume that  $R$  is the approximate inverse matrix of  $D_n - B_n$ , by (3) we have

$$\begin{aligned} & \|(I - \mathcal{P}_nN\mathcal{K}r\hat{u}^{N-1})^{-1}\|_{\mathcal{L}(X_n,X_n)} \\ &= \|(I - D_n^{-1}B_n)^{-1}\|_2 \\ &= \|(D_n - B_n)^{-1}D_n\|_2 \\ &= \|(D_n - B_n)^{-1}R^{-1}RD_n\|_2 \\ &= \|(I + R(D_n - B_n) - I)^{-1}RD_n\|_2 \\ &\leq \frac{\|RD_n\|_2}{1 - \|R(D_n - B_n) - I\|_2} = M_1. \end{aligned}$$

In this way, the constant  $C_1$  is given by Theorem2,

$$C_1 = \frac{1 + M_1K_1}{1 - (1 + M_1K_1)L_1}.$$

Secondly, we get the constant  $C_2$  with norm estimation. Namely, the residual of the operator equation (4) is necessary. The norm estimation is given so that

$$\begin{aligned} & \|F(\hat{u})\|_X \\ &= \|\hat{u} - \mathcal{K}r\hat{u}^N - \mathcal{K}f\|_X \end{aligned}$$

$$\begin{aligned}
&= \left\| \hat{u} - \mathcal{P}_n \mathcal{K} r \hat{u}^N - \mathcal{P}_n \mathcal{K} f_h - (\mathcal{K} - \mathcal{P}_n \mathcal{K}) r \hat{u}^N \right. \\
&\quad \left. - \mathcal{P}_n \mathcal{K} (f - f_h) - (\mathcal{K} - \mathcal{P}_n \mathcal{K}) f \right\|_X \\
&\leq \|Res\|_{X_n} + \|\mathcal{K} - \mathcal{P}_n \mathcal{K}\|_{\mathcal{L}(X,X)} \|r\|_{\infty} \|\hat{u}^N\|_{X_n} \\
&\quad + \|\mathcal{P}_n \mathcal{K}\|_{\mathcal{L}(X,X_n)} \|f - f_h\|_X + \|\mathcal{K} - \mathcal{P}_n \mathcal{K}\|_{\mathcal{L}(X,X)} \|f\|_X \\
&= C_2.
\end{aligned}$$

Here,  $\|Res\|_{X_n}$  is the residual of the finite linear system. The norm  $\|f - f_h\|_X$  can be estimated by the interpolation theory.  $\|\hat{u}\|_{\infty}$  is the maximum norm of  $\hat{u}$ . Then, we have the residual of the operator equation.

Finally,  $C_3$  is computable. We note that  $v, w \in B(\hat{u}, \delta)$ , these are bounded by

$$\|v\|_X \leq \|\hat{u}\|_{X_n} + \delta, \quad \|w\|_X \leq \|\hat{u}\|_{X_n} + \delta.$$

Therefore, we have

$$\begin{aligned}
&\|F'(v) - F'(w)\|_X \\
&= \|\mathcal{N} \mathcal{K} r (v^{N-2} + v^{N-3} w + \dots + v w^{N-3} + w^{N-2}) (v - w)\|_X \\
&\leq N(N-1) \|\mathcal{K}\|_{\mathcal{L}(X,X)} \|r\|_{\infty} (\|\hat{u}\|_{X_n} + \delta)^{N-2} \|v - w\|_X.
\end{aligned}$$

Accordingly,  $C_3$  follows that

$$C_3 = N(N-1) \|\mathcal{K}\|_{\mathcal{L}(X,X)} \|r\|_{\infty} (\|\hat{u}\|_{X_n} + \delta)^{N-2}.$$

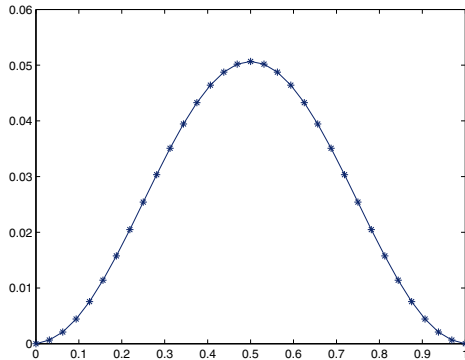
In this way, we have the constants  $C_1, C_2$  and  $C_3$ . The guaranteed error is bounded by the above verification method.

### 3. Numerical Example

For an application of our proposal method, we treated the following two-point boundary value problem:

$$\begin{cases} -u'' = u^3 - \cos 2\pi x & 0 < x < 1, \\ u(0) = u(1) = 0, \end{cases} \quad (7)$$

In this case, an approximate solution  $\hat{u}$  is led by FEM with Newton method. We choose that the finite element subspace  $S_h$  is one-dimensional piecewise hat functions. The shape of an approximate solution to (7) is shown as follows, where the divide number is 32.



For this approximate solution, we have

$$\alpha\omega < 1.8521 \times 10^{-8}.$$

Hence, there exists the exact solution of (7) in the ball centered at  $\hat{u}$  with radius:

$$\rho = 4.1557 \times 10^{-4}.$$

We have the guaranteed error:  $\|u - \hat{u}\|_X$  where  $X = L^2([0, 1])$ . Furthermore, increasing the divide number, we can improve the guaranteed error estimate.

Divide number	Guaranteed error estimate
8	6.649226221034610e-03
16	1.662284973637663e-03
32	4.155701919055114e-04
64	1.038924958465566e-04
128	2.597314309434374e-05
256	6.493342022539881e-06
512	1.623475629205407e-06
1024	4.063961081846812e-07
2048	1.026303844134345e-07

All computation is carried out on Mac OS X, Intel Core2 Duo 1.86GHz by using MATLAB 2009a with toolbox for verified computations, INTLAB[5].

### References

- [1] K. Atkinson and W. Han: Theoretical Numerical Analysis, Springer, 2001.
- [2] S.C. Brenner and L.R. Scott: The Mathematical Theory of Finite Element Methods, Springer, 2008.
- [3] M.T. Nakao and N. Yamamoto: Numerical Verification, Nihonhyouron-sya, 1998, (Japanese).
- [4] S. Oishi: Numerical Methods with Guaranteed Accuracy, Corona-sya, 2000, (Japanese).
- [5] S.M. Rump: INTLAB-INTerval LABoratory, a Matlab toolbox for verified computations, Hamburg University of Technology, "http://www.ti3.tu-harburg.de/rump/intlab/".
- [6] A. Takayasu, S. Oishi and T. Kubo: Numerical verification for solutions to nonlinear two-point boundary value problems with finite element method, proceedings of ITC-CSCC2009, 2009.
- [7] E. Zeidler: Nonlinear Functional Analysis and Its Applications, Part I Fixed Point Theorems, Springer-Verlag, 1986.