
ATTI ACCADEMIA NAZIONALE DEI LINCEI
CLASSE SCIENZE FISICHE MATEMATICHE NATURALI
RENDICONTI

SOLOMON G. MIKHLIN

**Some theorems on the stability of numerical
processes**

*Atti della Accademia Nazionale dei Lincei. Classe di Scienze Fisiche,
Matematiche e Naturali. Rendiconti, Serie 8, Vol. 72 (1982), n.2, p. 71–76.*

Accademia Nazionale dei Lincei

<http://www.bdim.eu/item?id=RLINA_1982_8_72_2_71_0>

L'utilizzo e la stampa di questo documento digitale è consentito liberamente per motivi di ricerca e studio. Non è consentito l'utilizzo dello stesso per motivi commerciali. Tutte le copie di questo documento devono riportare questo avvertimento.

*Articolo digitalizzato nel quadro del programma
bdim (Biblioteca Digitale Italiana di Matematica)
SIMAI & UMI*

<http://www.bdim.eu/>

Analisi numerica. — *Some theorems on the stability of numerical processes* (*). Nota (**) del Socio straniero SOLOMON G. MIKHLIN.

RIASSUNTO. — Nell'articolo si dimostrano alcuni teoremi sulla stabilità dei processi numerici di Ritz e della collocazione in rapporto agli errori di « distorsione ».

1. Let us consider a numerical process which consists in solving a sequence of independent equations

$$(1) \quad A_n x^{(n)} = f^{(n)}; \quad n = 1, 2, \dots$$

Here $x^{(n)} \in X_n, f^{(n)} \in Y_n$; A_n is an operator acting from X_n into Y_n ; X_n, Y_n are metric spaces. In this paper we only consider the case when X_n, Y_n are separable Banach spaces (in the sections 2-4—Hilbert spaces) and A_n are linear operators. Processes (1) arise, for example, when one uses the Ritz method (particularly, the finite elements method) for solving linear equations. In these cases the operators A_n and the right-hand terms $f^{(n)}$ are not given a priori. How it is natural, they are calculated with some errors. As a result we have to solve equations of a certain "distorted" sequence

$$(2) \quad (A_n + \Gamma_n) z^{(n)} = f^{(n)} + \delta^{(n)}$$

instead of sequence (1).

We say that the process (1) is stable, with respect to the distortions errors, in the sequence of pairs of spaces (X_n, Y_n) if there exist positive numbers p, q, r , such that the inequality $\|\Gamma_n\|_{X_n \rightarrow Y_n} \leq r$ involves the estimate

$$(3) \quad \|z^{(n)} - x^{(n)}\|_{X_n} \leq p \|\Gamma_n\|_{X_n \rightarrow Y_n} + q \|\delta^{(n)}\|_{Y_n}.$$

Some other definitions of stability are also possible.

It is demonstrated in [1] that the process (1) is stable, according to the above definition if and only if the conditions

$$(4) \quad \|A_n^{-1}\|_{Y_n \rightarrow X_n} \leq c_1, \|A_n^{-1} B_n x^{(n)}\|_{X_n} \leq c_2$$

are fulfilled; here c_1, c_2 do not depend on n , $x^{(n)}$ is the solution of (1) and B_n is an arbitrary operator with unit norm, acting from X_n into Y_n .

2. Let us consider the equation

$$(5) \quad Ax = f,$$

(*) Dedicated to Prof. G. Fichera on the occasion of his 60th birthday.

(**) Presentata nella seduta del 13 febbraio 1982.

where A is a positive definite [2] operator acting in a separable Hilbert space H ; we designate by H_A the energy space of the operator A , for the definition see [2]. We choose a sequence of finite-dimensional subspaces $H_A^{(n)} \subset H_A$; let this sequence be complete in H_A . We put $\dim H_A^{(n)} = N(n) = N$. Further let $(\varphi_{n1}, \varphi_{n2}, \dots, \varphi_{nN})$ be a basis in $H_A^{(n)}$. Following the Ritz method one constructs the approximate solution $x^{(n)}$ of (5) as an element of $H_A^{(n)}$

$$x^{(n)} = \sum_{k=1}^N a_k^{(n)} \varphi_{nk}$$

with coefficients $a_n^{(k)}$ satisfying the system of equations

$$(6) \quad M_n a^{(n)} = f^{(n)}.$$

Here M_n is the matrix of elements $[\varphi_{nk}, \varphi_{nj}]$; $a^{(n)}$ and $f^{(n)}$ are vectors in R_N with components $(a_1^{(n)}, a_2^{(n)}, \dots, a_N^{(n)})$ and $(f, \varphi_{n1}), (f, \varphi_{n2}), \dots, (f, \varphi_{nN})$ respectively. The indices j, k change in the limits $1 \leq j, k \leq N$; the square and round brackets designate the inner product in H_A and H respectively.

Remark. We obtain the classical Ritz method, if $\forall n, H_A^{(n)} \subset H_A^{(n+1)}$ [3]. The idea of using subspaces $H_A^{(n)} \not\subset H_A^{(n+1)}$ is due to Courant [4]; this idea contains the basis of the finite elements method.

Let A_n be the operator acting in R_N and generated by the matrix M_n . If $a^{(n)}$ and $f^{(n)}$ are treated as elements of R_N , then one can write the equation (6) in the form

$$(7) \quad A_n a^{(n)} = f^{(n)}.$$

It is demonstrated in [5] (see also [6]) that the numerical process (7) for the classical Ritz process is stable in the sequence (R_N, R_N) if and only if the least eigenvalue $\lambda_1^{(n)}$ of the matrix M_n is bounded below by a positive constant. The proof can be transferred without change on the case of non-expanding subspaces H_A .

3. We investigate now the stability of the Ritz process in the general case $\inf \lambda_1^{(n)} \geq 0$. We introduce two N -dimensional Hilbert spaces X_N and Y_N with the norms

$$(8) \quad \forall b \in R_N; \|b\|_{X_N} = \sqrt{\lambda_1^{(n)}} \|b\|_{R_N}, \|b\|_{Y_N} = \frac{1}{\sqrt{\lambda_1^{(n)}}} \|b\|_{R_N}.$$

Let us designate here by A_n the operator generated by the matrix M_n and acting from X_N into Y_N ; the vectors $a^{(n)}$ and $f^{(n)}$ are treated as elements of X_N and Y_N respectively.

THEOREM 1. *The process (7) is stable in the sequence (X_N, Y_N) .*

It is sufficient to prove that the inequalities (4) are satisfied.

Let $v^{(n)} \in H_A^{(n)}$, then

$$(9) \quad v^{(n)} = \sum_{k=1}^N b_k^{(n)} \varphi_{nk};$$

if we put $b^{(n)} = (b_1^{(n)}, b_2^{(n)}, \dots, b_N^{(n)})$, we obtain

$$(10) \quad \|v^{(n)}\|^2 = (M_n b^{(n)}, b^{(n)})_{R_N} \geq \lambda_1^{(n)} \|b^{(n)}\|_{R_N}^2 = \|b^{(n)}\|_{X_N}^2.$$

$\|\cdot\|$ designates the norm in H_A . Now

$$(11) \quad \|A_n^{-1}\|_{Y_N \rightarrow X_N} = \sup_{b \in R_N} \frac{\|A_n^{-1} b\|_{X_N}}{\|b\|_{Y_N}} = \lambda_1^{(n)} \sup_{b \in R_N} \frac{\|M_n^{-1} b\|_{R_N}}{\|b\|_{R_N}} = 1;$$

hence the first inequality (4) is proved.

The Ritz method converges in H_A for the equation (5), because A is positive definite [2]. Consequently, $\|x^{(n)}\| \leq c_3 = \text{const}$; according to (10), $\|a^{(n)}\|_{X_N} \leq c_3$. Now $\|A_n^{-1} B_n a^{(n)}\| \leq c_3$, and the second inequality (4) is also proved.

4. Formula (9) defines an operator Π_n which transforms any vector $b^{(n)} \in R_N$ in an element $v^{(n)} \in H_A^{(n)}$, so that $v^{(n)} = \Pi_n a^{(n)}$. The operator Π_n is invertible: $b^{(n)} = \Pi_n^{-1} v^{(n)}$; particularly, $a^{(n)} = \Pi_n^{-1} x^{(n)}$. Substituting this in (7), we obtain the numerical process giving the approximate solution $x^{(n)}$:

$$(12) \quad A_n \Pi_n^{-1} x^{(n)} = f^{(n)}.$$

THEOREM 2. *The numerical process (12) is stable in the sequence $(H_A^{(n)}, Y_N)$.*

We use the method of [7] in order to prove Theorem 2.

Let Γ_n and $\delta^{(n)}$ be the distortions of A_n and $f^{(n)}$ respectively, and let $c^{(n)}$ be the solution of the distorted equation

$$(13) \quad (A_n + \Gamma_n) c^{(n)} = f^{(n)} + \delta^{(n)}.$$

The distorted approximate Ritz solution is $z^{(n)} = \Pi_n c^{(n)}$, and

$$\begin{aligned} \|z^{(n)} - x^{(n)}\|^2 &= (M_n (c^{(n)} - a^{(n)}), c^{(n)} - a^{(n)})_{R_N} \leq \\ &\leq \|A_n (c^{(n)} - a^{(n)})\|_{Y_N} \cdot \|c^{(n)} - a^{(n)}\|_{X_N}. \end{aligned}$$

According to Theorem 1 there exist numbers $p, q, r > 0$ with the following property: if $\|\Gamma_n\|_{X_N \rightarrow Y_N} \leq r$, then

$$\|c^{(n)} - a^{(n)}\|_{X_N} \leq p \|\Gamma_n\|_{X_N \rightarrow Y_N} + q \|\delta^{(n)}\|_{Y_N}.$$

It follows from (7) and (13) that

$$(A_n + \Gamma_n)(c^{(n)} - a^{(n)}) = (I_n + \Gamma_n A_n^{-1}) A_n (c^{(n)} - a^{(n)}) = \delta^{(n)} - \Gamma_n a^{(n)},$$

where I_n is the identical operator in Y_n . Let r' be a number in the interval $(0, 1)$, and let $\|A_n^{-1}\| \cdot \|\Gamma_n\|_{X_N \rightarrow Y_N} \leq r'$. Then

$$\|(I_n + \Gamma_n A_n^{-1})^{-1}\| \leq (1 - r')^{-1}$$

and

$$\|A_n(c^{(n)} - a^{(n)})\| \leq \frac{1}{1 - r'} [c_3 \|\Gamma_n\|_{X_N \rightarrow Y_N} + \|\delta^{(n)}\|_{Y_N}].$$

Now obviously

$$(14) \quad \|z^{(n)} - x^{(n)}\| \leq p' \|\Gamma_n\|_{X_N \rightarrow Y_N} + q' \|\delta^{(n)}\|_{Y_N},$$

where p', q' are suitable constants. Theorem 2 is proved.

Remark. One can define the norms in X_N, Y_N as follows:

$$(15) \quad \forall b \in R_N; \|b\|_{X_N} = \gamma(n) \|b\|_{R_N}, \|b\|_{Y_N} = \frac{1}{\gamma(n)} \|b\|_{R_N}.$$

Here $\gamma(n)$ is any positive function of n , satisfying the inequality

$$\forall b \in R_N, \|\Pi_n b\| \geq C\gamma(n) \|b\|_{R_N}; C = \text{const},$$

or, what is the same,

$$(16) \quad \lambda_1^{(n)} \geq C\gamma^2(n).$$

In particular, it is sufficient that $\mu_1^{(n)} \geq C\gamma^2(n)$, where $\mu_1^{(n)}$ is the least eigenvalue of the matrix of inner products $(\varphi_{nk}, \varphi_{nj})_H; j, k = 1, 2, \dots, N$. Theorems 1 and 2 with their proofs still hold, only the relation (11) must be replaced by the inequality $\|A_n^{-1}\|_{Y_N \rightarrow X_N} \leq C^{-1}$, where C is the constant of (16).

The theorems on stability of the finite elements method given in [8] are particular cases of the Theorems 1 and 2. The function $\gamma(n)$ used in [8] is equal to $h^{m/2}$, where h is the step of the net and m is the dimension of the space of coordinates.

5. We consider now the problem of stability of the collocation method; this method was first formulated in [9]. The main points of the collocation method are the following. Let be given the problem of solving the equation (5) where A is a linear operator acting from a Banach space X into a Banach space Y , so that the domain $D(A)$ and the range $R(A)$ are dense in X and Y respectively. We suppose that Y consists only of functions which are continuous on a certain compact $K \subset R_m$. We choose a sequence of finite-dimensional

subspaces $X_n \subset D(A)$ which is complete in X and put $\dim X_n = N(n) = N$, $Y_n = AX_n$. If the inverse operator A^{-1} exists then $\dim Y_n = N$. Further, if $\{\varphi_{nk}\}$, $1 \leq k \leq N$, is a basis in X_n and $\psi_{nk} = A\varphi_{nk}$, then $\{\psi_{nk}\}$ is a basis in Y_n .

Let us choose some points $t_k^{(n)} \in K$, $1 \leq k \leq N$, the so-called collocation knots. One constructs the approximate solution of (5) as an element of X_n :

$$(17) \quad x^{(n)} = \sum_{k=1}^N a_k^{(n)} \varphi_{nk};$$

the coefficients $a_k^{(n)}$ are to be defined from the algebraic system

$$(18) \quad \sum_{k=1}^N a_k^{(n)} \psi_{kn}(t_j^{(n)}) = f(t_j^{(n)}); \quad 1 \leq j \leq N.$$

6. Let $t_j^{(n)}$ be the vertices of a certain parallelepipedal net. Further let $h_k^{(n)}$ be the length of the edge of the parallelepiped which is parallel to the k -th coordinate axis. Suppose

$$c_1 h_k^{(n)} \leq h_n \leq c_2 h_k^{(n)} \quad ; \quad c_1, c_2 = \text{const}, h_n \xrightarrow{n \rightarrow \infty} 0.$$

One can write down (18) in the form (6); the meaning of the notations is obvious. We consider $a^{(n)}$ as an element of R_N and $f^{(n)}$ as an element of the N -dimensional Hilbert space F_{Ns} with the norm $\|\cdot\|_{F_{Ns}} = h^{s/2} \|\cdot\|_{R_N}$.

Let A_n be the operator generated by the matrix M_n and acting from R_N into F_{Ns} . One can write the system (18) in the form (7).

Let $s_n^{(1)}$ designate the least singular number of the matrix M_n , i.e., the least eigenvalue of the non-negative matrix $M_n^* M_n$.

THEOREM 3. *If $\|a^{(n)}\| \leq c_3 = \text{const}$ and $s_1^{(n)} \geq c_4 h_n^{-s}$, where $c_3, c_4 = \text{const} > 0$, then the process (7) for the collocation method is stable in the sequence (R_N, F_{Ns}) . If $s_1^{(n)} \leq \gamma(n) h_n^{-s}$, $\gamma(n) \xrightarrow{n \rightarrow \infty} 0$, then the same process is unstable.*

Any numerical process of the kind (1) is stable if and only if the conditions (5) are fulfilled. It is easy to see that $\|A_n^{-1}\| = h_n^{-s/2} \|M_n^{-1}\|_{R_N \rightarrow R_N}$. The greatest singular number of M_n^{-1} is equal to $1/s_1^{(n)}$, hence $\|M_n^{-1}\|_{R_N \rightarrow R_N} = 1/\sqrt{s_1^{(n)}}$; consequently $\|A_n^{-1}\| = (h^s s_1^{(n)})^{-1/2}$. If $s_1^{(n)} \leq \gamma(n) h_n^{-s}$ then $\|A_n^{-1}\| \geq 1/\sqrt{\gamma(n)} \xrightarrow{n \rightarrow \infty} \infty$, and the process (7) is unstable. On the contrary, if $s_1^{(n)} \geq c_4 h_n^{-s}$ then $\|A_n^{-1}\| \leq 1/\sqrt{c_4}$ and the first condition (5) is satisfied. The second condition (5) is satisfied by assumption, and the collocation process is stable.

REFERENCES

- [1] MIKHLIN S. G. (1964) - *On the stability of some numerical processes*. «Dokl. Akad. Nauk SSSR», 157, N. 2, 271-273. English translation: «Soviet Math. Dokl.», 5, 931-933.
- [2] MIKHLIN S. G. (1957-1964) - *Variational methods in mathematical Physics*. Moscow. English translation: Pergamon Press, Oxford.
- [3] RITZ W. (1908) - *Über eine neue Methode zur Lösung gewisser Variations-probleme der mathematischen Physik*. Journ. f. d. reine und angewandte Mathematik », 135, H. 1.
- [4] COURANT R. (1943) - *Variational methods for the solution of problems of equilibrium and vibrations*. «Bull. Amer. Math. Soc. », 49, 1-23.
- [5] MIKHLIN S. G. (1960-1961) - *On the stability of the method of Ritz*. «Dokl. Akad. Nauk SSSR », 135, 16-19. English translation: «Soviet Math. Dokl. », 1, 1230-1233.
- [6] MIKHLIN S. G. (1971) - *The numerical performance of variational methods*, Moscow, 1966. English translation: Wolters-Noordhoff Publ., Groningen.
- [7] YASKOVA G. N., YAKOVLEV M. N. (1962) - *Some condition for the stability of the method of Petrov - Galerkin*. «Trudy Matem. in-ta im.». Steklova, 66, 182-189.
- [8] MIKHLIN S. G. (1979) - *Approximation on a rectangular grid*. Sijthoff & Noordhoff, Alphen aan den Rijn.
- [9] KANTOROVICH L. V. (1934) - *On a method for approximate solution of differential equations*. «Dokl. Akad. Nauk SSSR », 2, 532-536.