

On a Reasonable Choice of the Coordinate Functions in the Faedo-Galerkin Method

by Jan BOCHENEK

Introduction. In Mikhlin's monograph [7] conditions of a reasonable choice of coordinate functions for an approximate solution of the operator equation

$$(1) \quad Au = f,$$

were given with the use of the Ritz method. It is shown in [7] that on the choice of these coordinate functions, among other, the following properties depend:

1° convergence of the approximate solution and, possibly, convergence of the "residuum $Au_n - f$ " to zero,

2° rate of convergence of the approximating sequence.

In the papers [1], [2], [3] this result of Mikhlin is generalized and transferred to the Bubnov-Galerkin method for the equation of the type (1), when the operator A is not self-adjoint.

The purpose of this paper is to give a certain fashion of a reasonable choice of the coordinate functions in the method of Faedo-Galerkin for the approximate solution of the equation of the parabolic type (cf. [4]).

Introductory concepts. Let X be a real separable Hilbert space with the norm $|\cdot|$ and the inner product (\cdot, \cdot) and let A be a linear operator with domain $D(A) \subset X$. We assume that A is a self-adjoint positive definite operator i.e., the operator A satisfies the following condition

$$(2) \quad (Au, u) \geq p|u|^2, \quad u \in D(A), \quad p > 0.$$

Let us denote by $L^2(a, b; X)$ the linear space of the functions $f: [a, b] \rightarrow X$, such that

$$\int_a^b |f(t)|^2 dt < \infty,$$

where $|f(t)|$ denotes the norm of the element $f(t) \in X$, and $-\infty \leq a < b \leq +\infty$. The space $L^2(a, b; X)$ is a Hilbert space with the norm

$$(3) \quad \|f\| := \left\{ \int_a^b |f(t)|^2 dt \right\}^{1/2}$$

and with the inner product

$$(4) \quad \langle f, g \rangle := \int_a^b (f(t), g(t)) dt,$$

in particular,

$$L^2(X) := L^2(-\infty, +\infty; X).$$

Let us denote by $D^k(a, b; X)$, for arbitrary integer $k \geq 0$, the set of functions $f: [a, b] \rightarrow X$ which are of the class C^k in (a, b) such that $f, f', \dots, f^{(k)} \in L^2(a, b, X)$. In the set $D^k(a, b; X)$ we introduce the structure of a linear space and the norm by formula

$$(5) \quad \|f\|_{D^k(a, b; X)} := \left\{ \sum_{p=0}^k \int_a^b |f^{(p)}(t)|^2 dt \right\}^{1/2}.$$

The closure of the space $D^k(a, b; X)$ in the norm defined by (5) is denoted by $H^k(a, b; X)$.

Let $0 < T \leq +\infty$ be fixed. We shall denote by $H_0^k(0, T; X)$ the closure of the subspace $D_0^k(-\infty, T, X) := \{f \in D^k(-\infty, T; X) : f(t) = f'(t) = \dots = f^{(k)}(t) = 0 \text{ for } t < 0\}$ in the norm defined by (5).

In this paper the spaces $H^0(0, T; X) = L^2(0, T; X)$, $H^1(0, T; X)$ and $H_0^1(0, T; X)$ play a substantial part. Since all these spaces are subspaces of the space $L^2(0, T; X)$, we shall assume that the norm and the inner product of these spaces are defined by the formulas

$$(6) \quad \|f\| := \left\{ \int_0^T |f(t)|^2 dt \right\}^{1/2}$$

$$(7) \quad \langle f, g \rangle := \int_0^T (f(t), g(t)) dt.$$

Let P be a linear operator such that

$$(8) \quad P: H_0^1(0, T; X) \rightarrow H^0(0, T; X),$$

defined by the formula

$$(9) \quad Pu := \frac{du}{dt} + \mathcal{A}u \quad \text{for } u \in H_0^1(0, T; X),$$

where \mathcal{A} is a linear operator such that

$$\mathcal{A} := L^2(0, T; X) \rightarrow L^2(0, T; X)$$

defined by the formula

$$(10) \quad (\mathcal{A}u)(t) := Au(t) \text{ for } t \in [0, T] \text{ and } u \in L^2(0, T; X),$$

where A is the operator defined above. We shall prove the following

LEMMA 1. *The operator P defined by the formula (9) has the inverse bounded operator P^{-1} in the space $L^2(0, T; X)$.*

Proof. Let u be an element of $H_0^1(0, T; X)$. If $T < +\infty$, then $u: [0, T] \rightarrow X$ is a continuous function, therefore $u(0)$ and $u(T)$ have sense and $u(0) = 0$ (cf. [6], Ch. 1). By (2) and (4) we have

$$\begin{aligned} \langle Pu, u \rangle &= \int_0^T (u'(t), u(t)) dt + \int_0^T (Au(t), u(t)) dt \geq \frac{1}{2} \int_0^T \frac{d}{dt} |u(t)|^2 dt \\ &\quad + p \int_0^T |u(t)|^2 dt = \frac{1}{2} |u(T)|^2 + p \|u\|^2 \geq p \|u\|^2. \end{aligned}$$

If, however, $T = +\infty$, then for each $T_1 < +\infty$ we have

$$\int_0^{T_1} (u'(t) + Au(t), u(t)) dt \geq p \int_0^{T_1} |u(t)|^2 dt.$$

From this, letting $T_1 \rightarrow +\infty$, we get

$$\langle Pu, u \rangle \geq p \|u\|^2.$$

Therefore, for each $T \leq +\infty$, we get

$$(11) \quad \langle Pu, u \rangle \geq p \|u\|^2.$$

Putting $v = Pu$ in (11) we have

$$(12) \quad \|P^{-1}v\| \leq \frac{1}{p} \|v\|.$$

The proof is completed.

Apart from the operator A we now consider a linear self-adjoint, positive definite operator B such that $D(B) = D(A)$. We assume that the operator B possesses a discrete spectrum. We denote by $\{\lambda_n\}$ the denumerable sequence of eigenvalues of the operator B , and by $\{\varphi_n\}$ the corresponding orthonormal sequence of eigenvectors which form a complete system in the space X .

From the assumptions on the operators A and B it follows that AB^{-1} , BA^{-1} , $A^{-1}B$ and $B^{-1}A$ are bounded operators in X (cf. [7], Ch. I).

Let us denote by \mathcal{B} the linear operator such that

$$(13) \quad \mathcal{B} := L^2(0, T; X) \rightarrow L^2(0, T; X)$$

defined by the formula

$$(14) \quad (\mathcal{B}u)(t) := Bu(t) \text{ for } t \in [0, T] \text{ and } u \in L^2(0, T; X),$$

where B is the operator defined above.

It is easy to prove that $\mathcal{A}\mathcal{B}^{-1}$, $\mathcal{B}\mathcal{A}^{-1}$, $\mathcal{A}^{-1}\mathcal{B}$ and $\mathcal{B}^{-1}\mathcal{A}$ are bounded operators in the space $L^2(0, T; X)$.

LEMMA 2. The operator $P^{-1}\mathcal{B}$ is a bounded operator in the space $L^2(0, T; X)$.

Proof. First let us observe that for the operator P there exists the adjoint operator P^* , defined by the formula

$$(15) \quad P^*v = -\frac{dv}{dt} + \mathcal{A}v \quad \text{for } v \in D(P^*),$$

where

$$D(P^*) = \{v \in H^1(0, T; X) : v(T) = 0\} \quad \text{when } T < +\infty$$

and

$$D(P^*) = H^1(0, +\infty; X) \quad \text{when } T = +\infty.$$

Hence $D(P^*)$ is a dense subset of the space $H^1(0, T; X)$ and, for any $v \in D(P^*)$, we have

$$\begin{aligned} \langle P^*v, v \rangle &= -\int_0^T (v'(t), v(t)) dt + \int_0^T (\mathcal{A}v(t), v(t)) dt \geq -\frac{1}{2}|v(t)|^2 \Big|_0^T \\ &\quad + p \int_0^T |v(t)|^2 dt = \frac{1}{2}|v(0)|^2 + p\|v\|^2 \geq p\|v\|^2. \end{aligned}$$

Therefore the operator P^* has the inverse bounded operator $(P^*)^{-1}$ and

$$(16) \quad \|(P^*)^{-1}\| \leq \frac{1}{p}.$$

We shall prove that $R(P^*) = D((P^*)^{-1})$ is the whole space $L^2(0, T; X)$. To this purpose it is sufficient to prove that $R(P^*)$ is dense in $D(P^*)$. Let $z \in D(P^*)$ and $y \in R(P^*)$. We shall prove that if $\langle z, y \rangle = 0$ for any $y \in R(P^*)$, then $z = 0$. Indeed, if $y \in R(P^*)$, then $y = P^*x$, where $x \in D(P^*)$, whence $0 = \langle y, z \rangle = \langle P^*x, z \rangle = \langle x, 0 \rangle = \langle x, Pz \rangle$, i.e., $Pz = 0$. On the other hand, by (11) we have $p\|z\|^2 \leq \langle Pz, z \rangle = 0$ whence $z = 0$. It means that $R(P^*)$ is dense in $D(P^*)$. Since $D(P^*)$ is dense in $H^1(0, T; X)$ and $H^1(0, T; X)$ is dense in $L^2(0, T; X)$, then $R(P^*)$ is dense in $L^2(0, T; X)$. Since P^* is a closed operator, whence $R(P^*)$ is a closed and dense domain in $L^2(0, T; X)$. Therefore $R(P^*) = D((P^*)^{-1}) = L^2(0, T; X)$.

From the definition of the operator \mathcal{B} it follows that $D(P^*) \subset D(\mathcal{B})$. Therefore we can define $\mathcal{B}(P^*)^{-1}$. Since the operator $\mathcal{B}(P^*)^{-1}$ is closed and defined in the whole space $L^2(0, T; X)$, it is bounded (cf. [8], p. 557). It follows that the operator $P^{-1}\mathcal{B}$ is bounded as the adjoint operator to $\mathcal{B}(P^*)^{-1}$.

The Faedo-Galerkin method. The Faedo-Galerkin method is an approximate method which is used, among other things, to solve equations of the type

$$(17) \quad Pu = f,$$

where f is an element of the space $L^2(0, T; X)$, and P is the operator defined by the formula (9).

Let $\{x_n\}$ be an orthonormal and complete system of vectors of the space X . Let us denote

$$(18) \quad u_n(t) := \sum_{k=1}^n c_k(t) x_k \quad \text{for } t \in [0, T],$$

where c_k ($k = 1, \dots, n$) are real-valued functions defined on the interval $[0, T]$, such that for every $t \in (0, T)$ the following equations hold

$$(19) \quad (Pu_n(t), x_k) = (f(t), x_k), \quad k = 1, \dots, n,$$

where $Pu(t) := u'(t) + Au(t)$ for every $u \in H_0^1(0, T; X)$ and $t \in (0, T)$.

The system (19), as follows from the definition of the operator P , is a linear system of n ordinary differential equations of the first order with the unknown functions c_1, \dots, c_n . The system (19) is called the Faedo-Galerkin system for the equation (17). The sequence $\{u_n\}$ defined by (18), where the coefficients c_k ($k = 1, \dots, n$) satisfy (19), is called the approximate sequence for the equation (17) in the sense of Faedo-Galerkin. It is proved that for any orthonormal and complete system $\{x_n\}$ in X , the sequences $\{u_n(t)\}$ and $\{u_n'(t)\}$ are convergent to $u(t)$ and $u'(t)$, respectively, for every $t \in (0, T)$, where u is the unique solution of the equation (17) (cf. [9]). We shall prove the following

THEOREM 1. *If in the Faedo-Galerkin method for the equation (17) we take the sequence $\{\varphi_n\}$ of the eigenvectors of the operator B , defined above, then for every $t \in (0, T)$*

$$(20) \quad |Pu_n(t) - f(t)| \rightarrow 0 \quad \text{if } n \rightarrow +\infty.$$

Proof. Let us denote by M_n and L_n the closed linear subspaces spanned by $\{\varphi_1, \dots, \varphi_n\}$ and $\{\psi_1, \dots, \psi_n\}$, respectively, and by P_n and π_n the projectors of X onto M_n and L_n , respectively, where $\psi_k = AB^{-1}\varphi_k$, $k = 1, \dots, n$. By the definition of the operator P and by (18) the Faedo-Galerkin system (19) takes the form

$$(21) \quad c_j'(t) + \sum_{k=1}^n c_k(t) (A\varphi_k, \varphi_j) = (f(t), \varphi_j), \quad j = 1, \dots, n.$$

Let us denote

$$(22) \quad z_n(t) := Au_n(t)$$

and let us observe that

$$z_n(t) = Au_n(t) = \sum_{k=1}^n c_k(t) A\varphi_k = \sum_{k=1}^n c_k(t) \lambda_k AB^{-1}\varphi_k = \sum_{k=1}^n a_k(t) \psi_k,$$

where $a_k(t) = \lambda_k c_k(t)$, $\psi_k = AB^{-1}\varphi_k$, $k = 1, \dots, n$.

Now the system (21) may be written in the form

$$(23) \quad c_j'(t) + \sum_{k=1}^n a_k(t) (\psi_k, \varphi_j) = (f(t), \varphi_j), \quad j = 1, \dots, n,$$

or in another form

$$(24) \quad P_n u_n'(t) + P_n z_n(t) = P_n f(t), \quad t \in (0, T).$$

In the paper [2] it was proved that for every $v_n \in L_n$, we have the inequality

$$(25) \quad |v_n| \leq C \lambda_n^{\frac{1}{2}} |B^{-\frac{1}{2}} P_n v_n|, \quad C > 0.$$

Since $z_n(t) - \pi_n[f(t) - u'_n(t)] \in L_n$, by (25) we get

$$|z_n(t) - \pi_n[f(t) - u'_n(t)]| \leq C \lambda_n^{\frac{1}{2}} |B^{-\frac{1}{2}} P_n \{z_n(t) - \pi_n[f(t) - u'_n(t)]\}|.$$

From this, by (24), we have

$$(26) \quad |z_n(t) - \pi_n[f(t) - u'_n(t)]| \leq C \lambda_n^{\frac{1}{2}} |B^{-\frac{1}{2}} P_n \{[f(t) - u'_n(t)] - \pi_n[f(t) - u'_n(t)]\}|.$$

Owing to the paper [2], the operators $B^{-\alpha}$ and P_n are commutative for every $\alpha > 0$ and $n \in N$. From this and from inequality $|P_n| \leq 1$, by (26), we have

$$(27) \quad |z_n(t) - \pi_n[f(t) - u'_n(t)]| \leq C \lambda_n^{\frac{1}{2}} |B^{-\frac{1}{2}} \{[f(t) - \pi_n f(t)] - [u'_n(t) - \pi_n u'_n(t)]\}|.$$

Making the most of the boundedness of the operators AB^{-1} and $B^{-1}A$, the completeness of the sequence $\{\varphi_n\}$ and $\{\psi_n\}$ in the space X , the convergence of the sequence $\{u'_n(t)\}$ in the norm of the space X and reasoning analogously as in the proof of the Theorem 1 in the paper [2], we can get

$$(28) \quad \lim_{n \rightarrow \infty} \lambda_n^{\frac{1}{2}} |B^{-\frac{1}{2}} \{[f(t) - \pi_n f(t)] - [u'_n(t) - \pi_n u'_n(t)]\}| = 0$$

From the inequality (27) by (28) follows that

$$(29) \quad \lim_{n \rightarrow \infty} |z_n(t) - \pi_n[f(t) - u'_n(t)]| = 0.$$

Let us observe that

$$(30) \quad \begin{aligned} |u'_n(t) + z_n(t) - f(t)| &= |f(t) - u'_n(t) - z_n(t)| \\ &= |\{[f(t) - u'_n(t)] - \pi_n[f(t) - u'_n(t)]\} - \{z_n(t) - \pi_n[f(t) - u'_n(t)]\}| \\ &\leq |[f(t) - u'_n(t)] - \pi_n[f(t) - u'_n(t)]| + |z_n(t) - \pi_n[f(t) - u'_n(t)]|, \end{aligned}$$

and

$$(31) \quad \begin{aligned} |[f(t) - u'_n(t)] - \pi_n[f(t) - u'_n(t)]| &\leq |f(t) - \pi_n f(t)| + |u'_n(t) - u'(t)| \\ &\quad + |u'(t) - \pi_n u'(t)| + |\pi_n[u'_n(t) - u'(t)]| \leq |f(t) - \pi_n f(t)| \\ &\quad + 2|u'_n(t) - u'(t)| + |u'(t) - \pi_n u'(t)|, \end{aligned}$$

where $u'(t)$ is the limit of the sequence $\{u'_n(t)\}$ for every $t \in (0, T)$, in the norm of the space X .

From the inequalities (30) and (31) by (29) and the completeness of the sequence $\{\psi_n\}$, we get the equality

$$\lim_{n \rightarrow \infty} |u'_n(t) + z_n(t) - f(t)| = 0 \quad \text{for } t \in (0, T),$$

which is equivalent to (20) i.e., to the thesis of Theorem 1.

The rate of convergence of approximation in the sense of Faedo-Galerkin. Analogously as in the case of the equation (1), we shall prove the following

THEOREM 2. *If the operators $A, B, \mathcal{A}, \mathcal{B}$ and P are the operators defined in section 1, then*

$$(32) \quad \|u_n - u\| = o(\lambda_n^{-1}),$$

where $\{u_n\}$ is defined by (18), u is the unique solution of the equation (17), but $\{\lambda_n\}$ is the sequence of eigenvalues of the operator B .

Proof. Let us denote

$$(33) \quad \delta_n = Pu_n - f$$

and let us observe that for every $t \in (0, T)$, we have

$$\delta_n(t) = \sum_{k=1}^{\infty} (\delta_n(t), \varphi_k) \varphi_k.$$

By orthonormalization of $\{\varphi_n\}$ in the space X , we get

$$|\delta_n(t)|^2 = \sum_{k=1}^{\infty} |(\delta_n(t), \varphi_k)|^2.$$

On the other hand, by (19), for $x_k = \varphi_k$, $k = 1, \dots, n$, we get

$$(\delta_n(t), \varphi_k) = (Pu_n(t) - f(t), \varphi_k) = 0 \quad \text{for } k = 1, \dots, n.$$

It means that

$$(34) \quad |\delta_n(t)|^2 = \sum_{k=n+1}^{\infty} |(\delta_n(t), \varphi_k)|^2.$$

Now we shall estimate $|B^{-1}\delta_n(t)|$. We have

$$\begin{aligned} |B^{-1}\delta_n(t)|^2 &= \sum_{k=1}^{\infty} |(B^{-1}\delta_n(t), \varphi_k)|^2 = \sum_{k=1}^{\infty} |(\delta_n(t), B^{-1}\varphi_k)|^2 \\ &= \sum_{k=1}^{\infty} \frac{|(\delta_n(t), \varphi_k)|^2}{\lambda_k^2} = \sum_{k=n+1}^{\infty} \frac{|(\delta_n(t), \varphi_k)|^2}{\lambda_k^2} \leq \frac{|\delta_n(t)|^2}{\lambda_{n+1}^2}. \end{aligned}$$

Therefore

$$(35) \quad |B^{-1}\delta_n(t)| \leq \frac{|\delta_n(t)|}{\lambda_{n+1}}.$$

From the estimation (35), by the definition of the operator \mathcal{B} (cf. (13) and (14)), and by (6) we have

$$\|\mathcal{B}^{-1}\delta_n\|^2 = \int_0^T |(\mathcal{B}^{-1}\delta_n)(t)|^2 dt = \int_0^T |B^{-1}\delta_n(t)|^2 dt \leq \frac{1}{\lambda_{n+1}^2} \int_0^T |\delta_n(t)|^2 dt = \frac{\|\delta_n\|^2}{\lambda_{n+1}^2},$$

that is

$$\|\mathcal{B}^{-1}\delta_n\| \leq \frac{\|\delta_n\|}{\lambda_{n+1}}.$$

On the other hand, we have

$$(36) \quad \begin{aligned} \|u_n - u\| &= \|P^{-1}P(u_n - u)\| = \|P^{-1}(Pu_n - f)\| = \|P^{-1}\delta_n\| \\ &= \|P^{-1}\mathcal{B}\mathcal{B}^{-1}\delta_n\| \leq \|P^{-1}\mathcal{B}\| \|\mathcal{B}^{-1}\delta_n\| \leq \frac{\|P^{-1}\mathcal{B}\| \|\delta_n\|}{\lambda_{n+1}} \end{aligned}$$

In virtue of Lebesgue's theorem, we have

$$\lim_{n \rightarrow \infty} \|\delta_n\|^2 = \lim_{n \rightarrow \infty} \int_0^T |\delta_n(t)|^2 dt = \int_0^T \lim_{n \rightarrow \infty} |\delta_n(t)|^2 dt = 0.$$

Since by Lemma 2 $\|P^{-1}\mathcal{B}\| \leq C$, so (36) implies (32), that is the thesis of Theorem 2.

Theorem 2 gives an estimation of the rate of convergence of the sequence $\{u_n\}$ in the norm of the space $L^2(0, T; X)$. As it follows from the definition of the norm in the space $L^2(0, T; X)$, the convergences in the spaces $L^2(0, T; X)$ and X are not equivalent.

To prove the next theorem, we shall need the following

LEMMA Fatou (see [5]). *If $\{f_n\}$ is a sequence of nonnegative an integrable functions on the set $G \subset R^n$ which converges to a function f for almost every $x \in G$, then*

$$(37) \quad \int_G f(x) dx \leq \liminf_{n \rightarrow \infty} \int_G f_n(x) dx.$$

THEOREM 3. *Under the assumptions of the Theorem 2, the following estimation is true*

$$(38) \quad |u_n(t) - u(t)| = o(\lambda_n^{-1}),$$

for almost every $t \in (0, T)$.

Proof. Let us denote

$$(39) \quad f_n(t) := \lambda_n |u_n(t) - u(t)|, \quad t \in [0, T], \quad n \in N,$$

where λ_n , u_n and u have the same meaning as in Theorem 2. From the previous considerations it follows that for every $n \in N$

$$f_n: [0, T] \rightarrow R$$

is the continuous and nonnegative function in the interval $[0, T]$. From this, by (32), we have

$$(40) \quad \lim_{n \rightarrow \infty} \int_0^T f_n(t) dt = 0.$$

On the other hand because $u_n(t) \rightarrow u(t)$ for every $t \in [0, T]$, when $n \rightarrow \infty$, we have

$$\lim_{n \rightarrow \infty} f_n(t) = f(t) \quad \text{for every } t \in [0, T],$$

where f is any function defined in $[0, T]$. By (40) owing to Theorem 6 of [5], we get that the function f is integrable in $(0, T)$. From this, by Lemma Fatou, we have

$$(41) \quad \int_0^T f(t) dt \leq \liminf_{n \rightarrow \infty} \int_0^T f_n(t) dt = 0.$$

By (41), we get

$$(42) \quad \lim_{n \rightarrow \infty} f_n(t) = 0 \text{ for almost every } t \in (0, T).$$

The equality (42), by (39), is equivalent to the estimation (38), that is the thesis of Theorem 3.

References

- [1] J. Bochenek, *On the reasonable choice of the coordinate functions in the Bubnov-Galerkin method*, Comm. Math. 17 (1973), 9—16.
- [2] —, *Some remarks on the convergence of Ritz and Galerkin methods*, Comm. Math. 17 (1973), 17—27.
- [3] А. В. Джишкарини, *О методе Бубнова-Галеркина*, Жур. выч. Мат.-Физ. 7 (1967), Но. 2, 1398—1402.
- [4] S. Faedo, *Un nuovo metodo per l'analisi esistenziale e quantitativa dei problemi di propagazione*, Ann. Sc. Norm. Sup. Pisa (1949), 1—40.
- [5] S. Hartman, J. Mikusiński, *Teoria miary i calki Lebesgue'a*, PWN Warszawa 1957.
- [6] J. L. Lions, E. Magenes, *Problemes aux limites non homogènes et application*, vol. I., Dunod, Paris 1968.
- [7] С. Г. Михлин, *Численная реализация вариационных методов*, Москва 1966.
- [8] В. И. Смирнов, *Курс высшей математики*, т. 5, Москва 1959.
- [9] М. И. Вишик, *Задача Коши для уравнений с операторными коэффициентами, смешанная краевая задача для систем дифференциальных уравнений и приближенный метод их решения*, Матем. сб. 39 (81) вып. 1 (1956), 51—148.

Received November 2, 1981.