Hindawi Publishing Corporation Advances in Mathematical Physics Volume 2016, Article ID 7058017, 12 pages http://dx.doi.org/10.1155/2016/7058017



Research Article

The Ritz Method for Boundary Problems with Essential Conditions as Constraints

Vojin Jovanovic¹ and Sergiy Koshkin²

¹Systems, Implementation & Integration, Smith Bits, A Schlumberger Co., 1310 Rankin Road, Houston, TX 77032, USA ²Computer and Mathematical Sciences, University of Houston-Downtown, One Main Street #S705, Houston, TX 77002, USA

Correspondence should be addressed to Vojin Jovanovic; fractal97@hotmail.com

Received 22 November 2015; Accepted 17 February 2016

Academic Editor: Pavel Kurasov

Copyright © 2016 V. Jovanovic and S. Koshkin. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

We give an elementary derivation of an extension of the Ritz method to trial functions that do not satisfy essential boundary conditions. As in the Babuška-Brezzi approach boundary conditions are treated as variational constraints and Lagrange multipliers are used to remove them. However, we avoid the saddle point reformulation of the problem and therefore do not have to deal with the Babuška-Brezzi inf-sup condition. In higher dimensions boundary weights are used to approximate the boundary conditions, and the assumptions in our convergence proof are stated in terms of completeness of the trial functions and of the boundary weights. These assumptions are much more straightforward to verify than the Babuška-Brezzi condition. We also discuss limitations of the method and implementation issues that follow from our analysis and examine a number of examples, both analytic and numerical.

1. Introduction

In variational problems linear boundary conditions are often divided into essential (geometric) and natural (dynamic) [1, II.12] and [2, 4.4.7]. More generally, one calls the boundary conditions essential if they involve derivatives of order less than half of the order of the differential equation and natural otherwise [3, I.1.2]. In the standard exposition of the Ritz method the trial functions may violate the natural conditions but must satisfy all the essential ones [2, 4.4.7] and [4]. The reason is that the variational equations force the natural conditions on the trial solutions anyway, even if the trial functions themselves do not satisfy them.

But what if we wish to use trial functions that violate the essential conditions as well? For instance, in problems involving parametric asymptotics, the trial functions are preimposed with no regard for boundary conditions [5, 6] and in initial-boundary problems with time-dependent boundary conditions the (time independent) trial functions can not satisfy them in principle. One may also wish to use such violating trial functions because they are simpler; see [7] for other possible reasons. Thus, there is abundant motivation to generalize the Ritz method to trial functions that do not satisfy

the essential conditions. From a theoretical viewpoint this is a particular case of approximating solutions by nonconforming functions, the nonconformity here being at the boundary [8].

A natural idea is to treat the essential boundary conditions as variational constraints and to remove them as any other constraints using the Lagrange multipliers. Such approach is naively taken in some applied works at least since 1946 [9]; see also [10] where the authors explicitly cite the simplicity of the trial functions as a reason for using them. Babuška et al. [8, Sec.7] and [11] were the first to theoretically analyse the use of Lagrange multipliers in the context of nonconforming Finite Element Method (FEM); their work was generalized to more general trial functions by Brezzi [12] and [13, II.1]. Their analysis relies on the saddle point reformulation of the original variational problem and leads to the celebrated Babuška-Brezzi inf-sup condition that dictates a strict relation between choices of spaces for the trial functions and for the Lagrange multipliers; see also [13, 14] for applications in non-FEM context. This approach however is very far from the intuitive reasoning behind the naive application of the Lagrange multipliers in [9] or in [10]. As a result, the meaning of the Babuška-Brezzi condition

remains obscure, which explains why it remains relatively unfamiliar to nonexperts, and its verification is often quite involved mathematically [14]. The standard approach applies to problems with quadratic functionals [13, II.1] but produces very strong stability and approximation results.

It is not our purpose to match the technical sophistication of the Babuška-Brezzi approach, but to give a more elementary and natural derivation of the Ritz method with Lagrange multipliers, the Ritz-Lagrange method for short. It applies to general convex functionals and is much closer to the standard theory of the Ritz method [15-17]. We avoid the saddle point reformulation altogether and work directly with the original variational problem. A system of boundary weights is used to approximate the Lagrange multipliers, and our assumptions are stated in terms of completeness of the trial functions and of the boundary weights. While completeness is implicitly contained in the Babuška-Brezzi condition, its role is by no means obvious or well-known, see [18], and it has significant practical consequences. On the other hand, completeness is much more intuitive and straightforward to verify than an inf-sup condition. As a trade-off, we only prove convergence of the method rather than obtain an explicit error estimate, so the relation between spaces of the trial functions and of the boundary weights is less precise. But hopefully a more intuitive approach provides a better understanding of analytic and numerical issues involved.

This paper is largely inspired by observations in [18] on intriguing and counterintuitive effects that the lack of completeness has on convergence of trial solutions to a true solution. We illustrate these effects further in a number of examples and explain them in the general context of convex analysis. Our approach leads to a number of practically useful observations. In particular, extra care is needed compared to the usual Ritz method: the variational functional has to be more regular on a larger space, the trial functions have to be complete in this larger space as well, and the multipliers can not be eliminated from the approximating systems using the usual variational formulas because of convergence issues. In the higher dimensional problems one needs to balance the numbers of the boundary weights and the trial functions to obtain well-posed approximating problems. While we do not prescribe this balance precisely, which would involve an analog of the Babuška-Brezzi condition, we still obtain a practical rule of thumb that works well in examples.

In Section 2 we introduce the Ritz-Lagrange method using simple one-dimensional examples, where not only the exact solution, but even all trial solutions are computed analytically. We also introduce notation and terminology of convex analysis needed to analyse the method theoretically. The proof of convergence is obtained in this case by straightforward reduction to the classical Ritz method. Unfortunately, this direct approach does not carry over to higher dimensions, and we develop a suitable generalization in Section 3. Numerical applications to multidimensional problems follow in Section 4. The paper ends with Conclusions, where we summarize our findings and discuss Galerkin type generalizations. Technical proofs are collected in the Appendix.

2. Boundary Conditions as Variational Constraints

As a motivation, consider a boundary value problem for the second-order equation $u_{xx} = f$ on [0, L] with essential boundary conditions on both ends of the interval u(0) = u(L) = 0. We set for convenience $L = \pi$, and f = 1 to make everything explicitly computable. The exact solution is easily found to be $\overline{u} = (1/2)x^2 - (\pi/2)x$. The corresponding variational functional is

$$J(u) = \int_{0}^{L} \left(\frac{1}{2} (u_{x})^{2} + fu\right) dx, \tag{1}$$

and the boundary value problem is equivalent to minimizing it on functions satisfying the boundary conditions.

We select $\cos nx$, $n \ge 0$, as our trial functions; they obviously do not satisfy the boundary conditions. Taking N of them the trial solution is of the form $u^{(N)} = a_0/2 + \sum_{n=1}^N a_n \cos nx$ with unknown coefficients a_i (1/2 in front of a_0 is for agreement with the convention for the cosine series). Since our boundary conditions are essential, and our trial functions do not satisfy them, the Ritz method has to be modified in some way. Our approach is to treat the essential conditions as variational constraints and remove them using Lagrange multipliers. The Lagrange functional is $\mathfrak{L} = J + \lambda_0 u(0) + \lambda_\pi u(\pi)$, where λ_0 , λ_π are the Lagrange multipliers. Substituting, we find that

$$\mathfrak{L}\left(u^{(N)}\right) = \frac{\pi}{2}a_0 + \frac{\pi}{4}\sum_{n=1}^{N}n^2a_n^2 + \lambda_0\left(\frac{a_0}{2} + \sum_{n=1}^{N}a_n\right) + \lambda_\pi\left(\frac{a_0}{2} + \sum_{n=1}^{N}(-1)^na_n\right). \tag{2}$$

The variational equations are $\partial \Omega/\partial a_i = 0$, and adding the two boundary conditions we get the Ritz-Lagrange system. Solving these N+3 equations for N+3 unknowns $a_0, \ldots, a_N, \lambda_0$, and λ_{π} we get $\lambda_0 = \lambda_{\pi} = -\pi/2$, $a_n = (1+(-1)^n)/n^2$ for $n \ge 1$ and all N, while

$$a_0 = a_0^{(N)} := -4 \sum_{k=1}^{\lfloor N/2 \rfloor} \frac{1}{(2k)^2} = -\sum_{k=1}^{\lfloor N/2 \rfloor} \frac{1}{k^2} \xrightarrow[N \to \infty]{} -\frac{\pi^2}{6}, \quad (3)$$

where $\lfloor \cdot \rfloor$ is the floor function returning the largest integer not exceeding its argument. Thus, the trial solutions $u^{(N)}$ converge to the sum of the series:

$$u^{(N)}(x) \xrightarrow[N \to \infty]{} u^{(\infty)}(x) := -\frac{\pi^2}{12} + \sum_{n=1}^{\infty} \frac{1 + (-1)^n}{n^2} \cos nx.$$
 (4)

By extending the exact solution \overline{u} to an even function on $[-\pi, \pi]$ and expanding it into a cosine series with coefficients $a_n = (2/\pi) \int_0^\pi \overline{u}(x) \cos nx \, dx$ [19, 12.1], one finds that $u^{(\infty)}$ is exactly the cosine series of the exact solution $\overline{u} = (1/2)x^2 - (\pi/2)x$.

Variational problems typically come with a natural energy space, where convergence of solutions is considered [18]. On its energy space a variational functional J typically has two important properties: it is continuous and it is weakly coercive; that is, $J(u) \xrightarrow[\|u\| \to \infty]{} \infty$ [15, 6.2] and [16, III.10.2]. For a convex functional (we will only consider those) these

two properties are sufficient to prove convergence of the usual Ritz approximations in the energy norm [15, 6.2A]. For the functional from (1) the energy space is $W_2([0,\pi])$, which is the Hilbert space of functions square integrable with their first derivatives and vanishing at 0 and π , with the norm $\|u\|_{W_2^1} := (\|u\|_{L_2}^2 + \|u_x\|_{L_2}^2)^{1/2}$. This norm is stronger than the L_2 norm in the sense that any W_2^1 convergent sequence converges in L_2 , but not conversely.

For our purposes the concept of energy space is not quite suitable because it incorporates the essential boundary conditions, and our trial functions do not satisfy them. Instead we start with a reflexive Banach space $\mathscr U$ (the reader may assume it to be Hilbert without much loss) that has nothing to do with the boundary conditions, and a convex functional $J:\mathscr U\to\mathbb R$ on it. Next, we introduce the boundary operator, a linear map $\Gamma:\mathscr U\to\mathbb R^s$, that maps functions to their boundary values. The subspace $\mathscr U:=\{u\in\mathscr U\mid \Gamma u=0\}$ consists of functions that satisfy the boundary conditions. In our example we have $\mathscr U=W_2^1([0,\pi]), \Gamma u=(u(0),u(\pi))^T,$ and $\mathscr U=W_2^1([0,\pi])$. The following three assumptions turn $\mathscr U$ into a generalized analog of the energy space:

- (1) $J: \mathcal{U} \to \mathbb{R}$ is convex and continuous.
- (2) $J(u) \xrightarrow[\|u\| \to \infty]{} \infty$ on $\mathring{\mathcal{U}}$; that is, J is weakly coercive on $\mathring{\mathcal{U}}$.
- (3) $\Gamma: \mathcal{U} \to \mathbb{R}^s$ is linear and continuous.

This setup applies to homogeneous boundary conditions only. Nonhomogeneous conditions can be accommodated by selecting a function that satisfies them and switching to considering the differences with it. These differences solve the corresponding homogeneous problem, and all convergence issues can be reduced to them; see, for example, [20, 2.1].

We now turn to the trial functions. Recall that a system of elements in a Banach space is called complete if any element can be approximated by their linear combinations in the norm of the space. Let $\{\phi_i\}$ be a complete system in \mathcal{U} and let $\mathcal{U}^{(N)}$ denote the linear span of ϕ_1,\ldots,ϕ_N . The Ritz-Lagrange approach amounts to minimizing J on $\mathcal{U}^{(N)}$ subject to the boundary conditions, which is equivalent to minimizing it on $\mathring{\mathcal{U}}^{(N)} := \mathcal{U}^{(N)} \cap \mathring{\mathcal{U}}$. Although by assumption about completeness all elements in \mathcal{U} can be approximated by linear combinations of ϕ_i , it is not a priori clear that functions from $\mathring{\mathcal{U}}$ can be approximated by linear combinations of ϕ_i that are themselves in $\mathring{\mathcal{U}}$. The next lemma proved in the Appendix assures us that this is the case.

Lemma 1. For any complete system of elements in \mathcal{U} there exists a system of their finite linear combinations belonging to $\mathring{\mathcal{U}}$, which is complete in $\mathring{\mathcal{U}}$.

This lemma effectively reduces the Ritz-Lagrange method to the traditional Ritz method. Indeed, if $\{\widetilde{\phi}_i\}$ is the complete system in $\mathring{\mathcal{U}}$ produced by Lemma 1, then applying the Ritz method with $\widetilde{\phi}_i$ as the trial functions amounts to minimizing J on $\mathring{\mathcal{U}}^{(N)}$. In other words, the Ritz-Lagrange method with $\{\phi_i\}$ produces the same $u^{(N)}$ (up to reindexing) as the Ritz method with $\{\widetilde{\phi}_i\}$. This allows us to use well-known results on convergence of the Ritz method [15, 6.2A], [16, IV.12.4], and [21, 42.5] to prove convergence of its Ritz-Lagrange generalization. Producing the Ritz system involves differentiating the functional, so in addition to convexity and continuity we also have to assume its Gateaux differentiability [16, I.2.1], and [15, 3.2].

Theorem 2. Let $J: \mathcal{U} \to \mathbb{R}$ be a convex continuous Gateaux differentiable functional, let $\Gamma: \mathcal{U} \to \mathbb{R}^s$ be a bounded linear operator, and let $\{\phi_i\}$ be a complete system in \mathcal{U} . If J is weakly coercive on $\mathring{\mathcal{U}}$, then it has a minimizer \overline{u} on it, as well as minimizers $u^{(N)}$ on all $\mathring{\mathcal{U}}^{(N)}$, and there exists a subsequence N_k such that $u^{(N_k)} \xrightarrow[k \to \infty]{w} \overline{u}$, where $\xrightarrow[k \to \infty]{w}$ denotes weak convergence in \mathscr{U} . Moreover, if J is strictly convex on $\mathring{\mathcal{U}}$, then \overline{u} , $u^{(N)}$ are unique and $u^{(N)} \xrightarrow[k \to \infty]{w} \overline{u}$. In both cases the values of J converge to its minimum on $\mathring{\mathscr{U}}$.

The proof is fairly standard, but we outline it in the Appendix for the convenience of the reader. For general convex functionals only weak convergence can be expected; we discuss a stronger convexity assumption that guarantees convergence by norm in the next section. Note also that weak convergence in W_2^1 implies convergence by norm in L_2 due to Sobolev embedding theorems [20, I.6] and [22, 1.8].

Theorem 2 mostly justifies the Ritz-Lagrange method used to solve our example. The required properties of $J(u) = \int_0^L (1/2)(u_x)^2 + fu \, dx$ and the boundary operator $\Gamma(u) = (u(0), u(L))^T$ are easily verified, except for the strict convexity, which follows from the Poincaré-Friedrichs inequality [20, I.6]. The completeness is a trickier issue. Completeness in $L_2([0,\pi])$ follows from the standard theorems on cosine series [19, Ch.12], but we need a stronger form of completeness in W_2^1 . Fortunately, W_2^1 completeness of cosines reduces to the L_2 completeness of sines. The minimality below means that the system becomes incomplete after deleting any function.

Lemma 3. The system $\cos nx$, $n \ge 0$, is complete and minimal in $W_2^1([0,\pi])$.

The proof is given in the Appendix.

As emphasized in [18] completeness is not a mere technicality in this context; it imposes a practical restriction on the

choice of trial functions. To underscore the point, consider the biharmonic equation $u_{xxxx} = f$ with the boundary conditions $u(0) = u(\pi) = u_{xx}(0) = u_{xx}(\pi) = 0$. For f = 1 the exact solution is $\overline{u} = x^4/24 - \pi(x^3/12) + \pi^3(x/24)$. If the cosine system is used again, proceeding as above we find that

$$u^{(N)}(x) \xrightarrow[N \to \infty]{} u^{(\infty)}(x) := \frac{\pi^4}{720} - \sum_{n=1}^{\infty} \frac{1 + (-1)^n}{n^4} \cos nx = \frac{x^4}{24} - \pi \frac{x^3}{12}$$

$$+ \pi^2 \frac{x^2}{24}.$$
(5)

This time the trial solutions do not converge to the exact one; in fact $\overline{u}-u^{(\infty)}=\pi^3(x/24)-\pi^2(x^2/24)$. This is because *cosines are incomplete in* $W_2^2([0,\pi])$. As observed in [18], the second derivatives $0, -\cos x, -4\cos 2x, \ldots, -n^2\cos nx, \ldots$ do not include a constant and therefore can not approximate the second derivative of x^2 in L_2 . But then cosines can not approximate x^2 in W_2^2 since its norm incorporates the L_2 norm for second derivatives. In [18] the authors add x^2 to the cosine system, but as we just saw the limit difference is a linear combination of x^2 and x, so at least x also needs to be added. We show in the Appendix that nothing else is needed.

Lemma 4. The system x, x^2 , $\cos nx$, $n \ge 0$, is complete and minimal in $W_2^2([0,\pi])$.

Note that verifying completeness in a correct space can not be avoided even if one uses the usual Ritz method with trial functions satisfying the essential boundary conditions. The second example demonstrates that solutions can not always be approximated in the space where the trial functions happen to be complete. Only completeness in the norm dictated by the variational functional counts. Completeness of trial functions in a norm weaker than the energy norm does not simply weaken the convergence to the exact solution; trial solutions may not converge to it at all.

After adding x and x^2 as the trial functions the trial solutions become $u^{(N)} = b_1 x + b_2 x^2 + a_0/2 + \sum_{n=1}^{N} a_n \cos nx$ giving the Lagrange functional

$$\mathfrak{L}\left(u^{(N)}\right) = -\frac{\pi^2}{2}b_1 - \frac{\pi^3}{3}b_2 - \frac{\pi}{2}a_0 + 2\pi b_2^2 + \frac{\pi}{4}\sum_{n=1}^N n^4 a_n^2 + \lambda_0 \left(\frac{a_0}{2} + \sum_{n=1}^N a_n\right) + \lambda_\pi \left(\pi b_1 + \pi^2 b_2 + \frac{a_0}{2} + \sum_{n=1}^N (-1)^n a_n\right).$$
(6)

The Ritz-Lagrange system now has two extra variables b_1 , b_2 and two extra equations. Solving it we find that

$$u^{(\infty)}(x) = \frac{\pi^3}{24}x - \frac{\pi^2}{24}x^2 + \frac{\pi^4}{720} - \sum_{n=1}^{\infty} \frac{1 + (-1)^n}{n^4} \cos nx$$

$$= \overline{u}(x)$$
(7)

matches the exact solution as expected.

3. The Ritz-Lagrange Method in Higher Dimensions

The Ritz-Lagrange method described in Section 2 can not be applied to multidimensional boundary value problems. In this section we will develop a suitable generalization and prove that it works. The main distinction is that the boundary operator $\Gamma: \mathscr{U} \to \mathscr{V}$ no longer maps into a finite-dimensional space. Indeed, in dimensions two and higher the boundary values are not arrays of numbers but functions on the boundary forming an infinite-dimensional space \mathscr{V} . The induction proof of the key Lemma 1 no longer works and its claim itself is false. It is easy to find complete systems of functions with no (finite) nontrivial linear combinations satisfying the boundary conditions. If we are committed to using arbitrary complete systems of trial functions we must find a way to form their linear combinations that satisfy essential boundary conditions "approximately."

To this end we will use a complete system $\{\psi_j\}$ of linear functionals on the Banach space \mathcal{V} , that is, elements of the dual space \mathcal{V}^* (as with \mathcal{U} , the reader may assume that \mathcal{V} is a Hilbert space, in which case $\mathcal{V}^* = \mathcal{V}$). If $\langle \psi_j, \Gamma u \rangle = 0$ for all j, then $\Gamma u = 0$ and $u \in \mathcal{U}$, so we can think of operators $\Gamma_s(u) = (\langle \psi_1, \Gamma u \rangle, \dots, \langle \psi_s, \Gamma u \rangle)^T$ as approximations to Γ and the corresponding spaces $\mathcal{U}_s := \{u \in \mathcal{U} \mid \Gamma_s(u) = 0\}$ as approximations to \mathcal{U} . Assuming Γ is continuous Γ_s also will be and we can apply Theorem 2 with Γ_s in place of Γ for each S. This gives us a sequence of approximations $U_s^{(N)}$ converging to an exact minimizer \overline{u}_s of J on \mathcal{U}_s . The remaining question is whether we can count on \overline{u}_s to approximate the overall minimizer \overline{u} of J on \mathcal{U} . Before proceeding let us describe the approximating procedure that our approach suggests.

Multidimensional Ritz-Lagrange Method. To minimize a functional $J: \mathcal{U} \to \mathbb{R}$ subject to essential boundary conditions $\Gamma u = 0$ with $\Gamma: \mathcal{U} \to \mathcal{V}$ select internal trial functions $\phi_1, \ldots, \phi_N \in \mathcal{U}$ and boundary weight functions $\psi_1, \ldots, \psi_s \in \mathcal{V}^*$ with $N \gg s$. A Ritz-Lagrange trial solution $u_s^{(N)} = \sum_{i=1}^N a_i \phi_i$ is obtained by solving the system of N+s equations with N+s unknowns $a_1, \ldots, a_N, \lambda_1, \ldots, \lambda_s$ consisting of N internal equations $\partial \Omega/\partial a_i = 0$ and s boundary equations $\langle \psi_j, \Gamma u \rangle = 0$, where $\Omega := J(u_s^{(N)}) + \langle \lambda^{(s)}, \Gamma u_s^{(N)} \rangle$ is the Lagrange functional and $\lambda^{(s)} := \lambda_1 \psi_1 + \cdots + \lambda_s \psi_s$ is the Lagrange multiplier.

A justification of our approach is based on Theorem 6. The reader not interested in justification may skip the rest of this section and look at applications in the next one. Even in simplest cases we can not expect \overline{u}_s to converge to \overline{u} in the same generality as in Theorem 2. The root cause is that the minimizer in $\mathring{\mathscr{U}}$ is approximated by elements outside of $\mathring{\mathscr{U}}$, which is why we need J to be well-behaved on the entire \mathscr{U} , not just $\mathring{\mathscr{U}}$, something avoidable if ϕ_i do satisfy the boundary conditions. In particular, the values $J(\overline{u}_s)$ are potentially smaller than $J(\overline{u})$ because they are obtained by minimizing J on larger subspaces $\mathring{\mathscr{U}}_s \supset \mathring{\mathscr{U}}$. As a consequence, standard properties of convex functionals, which we relied on in

Theorem 2, no longer guarantee convergence of $J(\overline{u}_s)$ to $J(\overline{u})$ if \overline{u}_s converges only weakly (in technical terms, convex functionals are weakly lower semicontinuous, but not necessarily weakly upper semicontinuous [16, III.8.5] and [17, 41.2]).

To make our proof work we need to assume a stronger form of convexity of J. For Gateaux differentiable functionals J convexity is equivalent to monotonicity of their derivatives; that is, $\langle J'(u) - J'(v), u - v \rangle \geq 0$ for all u, v [16, II.5.3]. This is a generalization of a familiar fact that convex functions have monotone derivatives. We will need a form of uniform monotonicity for J', compare [16, VI.18.6] and [21, 25.3]; namely,

$$\langle J'(u) - J'(v), u - v \rangle \ge c (\|u - v\|),$$
 (8)

where c(t) is a continuous monotone increasing function with c(t) = 0. The point is that if $c(\|u_s\|) \to 0$ then $\|u_s\| \to 0$ and hence $u_s \to 0$ by norm. The next Lemma answers in the affirmative the question about convergence of intermediate minimizers \overline{u}_s under the uniform monotonicity assumption.

Lemma 5. Let $J: \mathcal{U} \to \mathbb{R}$ be a convex Gateaux differentiable functional, and let $\Gamma: \mathcal{U} \to \mathcal{V}$ be a bounded linear map. Let $\{\psi_j\}$ be a complete system in \mathcal{V}^* and set $\mathring{\mathcal{U}}_s := \{u \in \mathcal{U} \mid \langle \psi_j, \Gamma u \rangle = 0, \ 1 \leq j \leq s\}$. If J is weakly coercive on some $\mathring{\mathcal{U}}_{s_0}$, and J' is uniformly monotone on it, then J has a minimizer \overline{u} on $\mathring{\mathcal{U}}_s$ as well as minimizers \overline{u}_s on all $\mathring{\mathcal{U}}_s$ with $s \geq s_0$, and $\overline{u}_s \xrightarrow{s \to \infty} \overline{u}$.

Uniform monotonicity also allows us to improve on Theorem 2 by replacing weak limits with strong limits leading to our main result.

Theorem 6. Let $J: \mathcal{U} \to \mathbb{R}$ be a convex continuous Gateaux differentiable functional, and let $\Gamma: \mathcal{U} \to \mathcal{V}$ be a bounded linear map. Let $\{\phi_i\}$ be a complete system in \mathcal{U} , and let $\{\psi_j\}$ be a complete system in \mathcal{V}^* . Denote by $\mathcal{U}^{(N)}$ the linear span of ϕ_1, \ldots, ϕ_N , and set $\mathring{\mathcal{U}} = \{u \in \mathcal{U} \mid \Gamma u = 0\}, \mathring{\mathcal{U}}_s = \{u \in \mathcal{U} \mid \langle \psi_j, \Gamma u \rangle = 0, \ 1 \leq j \leq s\}$, and $\mathring{\mathcal{U}}_s^{(N)} = \mathcal{U}^{(N)} \cap \mathring{\mathcal{U}}_s$. Suppose that J is weakly coercive on some $\mathring{\mathcal{U}}_s$ and that J' is uniformly monotone on it. Then for $s \geq s_0$ J has unique minimizers \overline{u} , $u_s^{(N)}$ on $\mathring{\mathcal{U}}, \mathring{\mathcal{U}}_s^{(N)}$, respectively, and $\lim_{s \to \infty} \lim_{N \to \infty} u_s^{(N)} = \overline{u}$ in the norm of \mathcal{U} , while $\lim_{s \to \infty} \lim_{N \to \infty} J(u_s^{(N)}) = J(\overline{u})$.

In examples it is typical that J does not satisfy (8) on the entire space $\mathscr U$ but does satisfy it on subspaces much larger than $\mathring{\mathscr U}_s$. Let us discuss the case of quadratic functionals J(u)=(1/2)B(u,u)+l(u) in more detail. By direct calculation $J'(v)=B(v,\cdot)+l$, therefore,

$$\langle J'(u) - J'(v), u - v \rangle = B(u, u - v) - B(v, u - v)$$

= $B(u - v, u - v)$. (9)

We need $B(u, u) \ge \varepsilon ||u||^2$ with $\varepsilon > 0$ to satisfy (8); that is, we need B to be strictly positive definite. The multidimensional

analog of the functional from our first example is $J(u) = \int_{\Omega} (1/2)(\nabla u)^2 + fu \, dx$, where Ω is a domain with smooth boundary. It follows from the Poincaré-Friedrichs inequality [20, I.6] and [23] that $B(u,u) = \int_{\Omega} (1/2)(\nabla u)^2 dx$ is strictly positive definite on $W_2^1(\Omega)$, but it certainly is not on the entire $W_2^1(\Omega)$ since B(u,u) = 0 for any u = const. Nevertheless, it still follows from the calculus of variations that B satisfies (8) on any subspace complementary to the subspace of constants; see, for example, [24, VI.1]. Similar considerations apply to other quadratic forms related to the strongly elliptic equations like the biharmonic equation. They are usually strictly positive definite on complements to finite-dimensional subspaces that they annihilate [23, 22.11]. Such conditions are sometimes called Ker(Γ)-ellipticity [7, 12].

It should be emphasized that Theorem 6 does not imply that the double sequence $u_s^{(N)}$ converges to \overline{u} ; that is, the repeated limit in it can not be replaced by the double limit. In fact, suppose N < s and let the functionals $\Gamma^* \psi_1, \dots, \Gamma^* \psi_s$, where $\Gamma^*: \mathcal{V}^* \to \mathcal{U}^*$ is the dual of Γ , be linearly independent. Then the s boundary equations $\langle \psi_i, \Gamma u \rangle = \langle \Gamma^* \psi_i, \Gamma u \rangle$ $u\rangle = 0$ alone are enough to force $u_s^{(N)} = 0$ no matter how large N and s are. In practice, this means that one should always take many more internal trial functions than the boundary *ones*; hence the prescription $N \gg s$. This way for large N the approximation $u_s^{(N)}$ will be close to \overline{u}_s , while \overline{u}_s in turn will be close to \overline{u} if s itself is large enough. The readers familiar with non-conforming Finite Element methods will recognize this as a reflection of the Ladyzhenskaya-Babuška-Brezzi type condition. One of its consequences is that the mesh size on the boundary has to be larger than in the interior, yielding a smaller number of the boundary elements [7, 11]. Modern approach can be found in [25] for finite elements, and in [14] for Galerkin approximations.

As in one-dimensional examples one will have to verify completeness of trial functions in the appropriate space; the same applies to the boundary weights as well. One has to be extra careful with functionals involving higher-order derivatives because the values of function and their derivatives have to be approximated simultaneously. Natural spaces to use are $W_p^k(\Omega)$, the spaces of functions with integrable pth powers along with all of their derivatives up to order k. A generalization of the Weierstrass theorem implies that polynomials form a complete system in $W_p^k(\Omega)$ for any $p \geq 1$, any bounded domain Ω , and any positive integer k (in fact, polynomials are even uniformly complete [24, II.4.3]). However, polynomials may not always be convenient in a particular problem. The following lemma can be useful in finding other complete systems.

Lemma 7. Let $\{\phi_i\}$ and $\{\widetilde{\phi}_j\}$ be complete systems in $W_p^k(\Omega)$ and $W_p^k(\widetilde{\Omega})$, respectively, where Ω and $\widetilde{\Omega}$ are some bounded domains. Then the system $\{\phi_i\widetilde{\phi}_i\}$ is complete in $W_p^k(\Omega\times\widetilde{\Omega})$.

If one starts from one-dimensional systems the lemma will only produce complete systems in box-like domains $[a_1, b_1] \times \cdots \times [a_m, b_m]$. However, any system of functions

complete on a domain will be complete on any of its subdomains, so for an arbitrary domain one can always use a system complete on the smallest box that contains it. A more targeted choice is to take eigenfunctions of an operator on the same domain that is simpler than the one involved but is somewhat similar to it. Various spectral theorems often ensure completeness of eigenfunctions in suitable Sobolev spaces [23, 22.11a].

4. Multidimensional Examples

In this section we illustrate the multidimensional Ritz-Lagrange method developed in Section 3 by applying it to some typical problems. Since calculations by hand quickly become intractable we performed them using a computer algebra system.

Consider a boundary value problem for the Laplace equation $\nabla^2 u = f$ in Ω , where Ω is the unit disk, with the boundary condition u = 0 on $\partial\Omega$ and $f = \cos(\sqrt{x^2 + y^2})$. This equation describes the transverse deflection of a membrane fixed everywhere at the boundary and subjected to pressure given by f [24, IV.10.3]. The profile of f was chosen so that the problem has an analytic solution which is not a polynomial. Specifically, one can represent the exact solution as a rapidly convergent series

$$\overline{u}(r) = \gamma + \cos(1) - \text{Ci}(1) - \cos(r) + \sum_{i=1}^{\infty} \frac{(-1)^i}{2i(2i)!} r^{2i},$$
 (10)

where $r = \sqrt{x^2 + y^2}$, $\gamma := \lim_{n \to \infty} (\sum_{k=1}^n (1/k) - \ln n)$ is the Euler-Mascheroni constant, and $Ci(x) := -\int_x^\infty (\cos t/t) dt$ is the cosine integral.

To solve the problem we use the multidimensional Ritz-Lagrange method. A variational formulation is to minimize the functional $J(u) = \int_{\Omega} (1/2)(\nabla u)^2 + fu \, dx \, dy$, which gives the total potential energy of the membrane, subject to the boundary condition. In the notation of Section 3 we take $\mathcal{U} = W_2^1(\Omega)$ with Γ being the restriction of u to the boundary $\partial\Omega$. Moreover, Γ is continuous if we take $\mathscr{V}=L_2(\partial\Omega)$. Our internal trial functions are the monomials, which obviously do not satisfy the boundary condition, and the trial solution is $u^{(N)} = \sum_{i=1}^{N} \sum_{j=1}^{N} a_{ij} x^{i-1} y^{j-1}$. Note that N of Section 3 will be N^2 here because of double indexing. As the boundary weight functions we choose the piecewise linear ones on uniform partitions of $\partial\Omega$. Unlike some circle specific choices, for example, the trigonometric functions, such weights can be used on a wide variety of boundaries. Instead of using a single indexed system ψ_k it is convenient to split it into the constant terms g_k and the linear terms h_k . If the boundary is partitioned into s segments, we have

$$g_{k}(\theta) \coloneqq \begin{cases} 1, & \frac{2\pi k}{s} \le \theta \le \frac{2\pi (k+1)}{s}, \\ 0, & \text{otherwise}, \end{cases}$$

$$h_{k}(\theta) \coloneqq \begin{cases} \theta, & \frac{2\pi k}{s} \le \theta \le \frac{2\pi (k+1)}{s}, \\ 0, & \text{otherwise}. \end{cases}$$
(11)

Table 1: Central (x = 0; y = 0) and boundary (x = 0; y = 1) error relative to the maximum deflection of the membrane of unit radius.

| | Central error% | Boundary error% |
|--------------|----------------|-----------------|
| N = 3, s = 2 | -4.04 | -1.33 |
| N = 4, s = 3 | -4.05 | -0.1 |
| N = 5, s = 4 | -0.04 | 0.03 |
| | | |

Therefore, the number of boundary weight functions, denoted by s in Section 3, will be 2s here. The Lagrange multiplier has the form $\lambda^{(s)} = \sum_{k=0}^{s-1} (c_k g_k(\theta) + d_k h_k(\theta))$, and the Lagrange functional is $\mathfrak{L}(u^{(N)}) = J(u^{(N)}) + \int_{\partial\Omega} \lambda^{(s)} u^{(N)} d\sigma$. The unknown coefficients a_{ij}, c_k , and d_k are determined from the system of N^2 internal $\partial \mathfrak{L}/\partial a_{ij} = 0$ and 2s boundary $\int_{\partial\Omega} g_k u^{(N)} d\sigma = \int_{\partial\Omega} h_k u^{(N)} d\sigma = 0$ equations. The relative errors of the Ritz-Lagrange solutions versus

The relative errors of the Ritz-Lagrange solutions versus the exact solution \overline{u} (10) are shown in Table 1 as the percentages of the maximum deflection at x=0 and y=0. They are quite small considering that one has to determine N^2+2s coefficients in each case. Note that we always keep $N^2>2s$ as recommended in the description of the method to ensure that the system matrices have full rank and are invertible.

Our next example involves a fourth-order equation. Consider the problem of bending a uniformly loaded, simply supported on all sides (SS-SS-SS), isotropic, square plate of constant thickness, unit stiffness, and unit edge length. Simply supported means that u=0 on $\partial\Omega$. We do not need to list the natural boundary conditions since a variational formulation incorporates them automatically. The variational functional giving the potential energy of the plate is as follows [24, IV.10.3]:

$$J(u) = \iint_0^1 \left(\frac{1}{2} \left(\left(\frac{\partial^2 u}{\partial x^2} \right)^2 + \left(\frac{\partial^2 u}{\partial y^2} \right)^2 + 2v \frac{\partial^2 u}{\partial x^2} \frac{\partial^2 u}{\partial y^2} \right) \right) dx dy,$$

$$(12)$$

where u is the displacement of the plate, v is the Poisson ratio, and f is a distributed load.

The Euler-Lagrange equation induced by (12) is the biharmonic equation $\nabla^2 \nabla^2 u = f$; the terms multiplied by the Poisson ratio form a divergence and only affect the natural boundary conditions. For the trial functions we choose the products of cosines $X_i(x) = \cos((i-1)\pi x)$ and $Y_i(y) = \cos((i-1)\pi y)$, so that the trial solution is $u^{(N)} = \sum_{i=1}^N \sum_{j=1}^N a_{ij} X_{i-1}(x) Y_{j-1}(y)$. Obviously, the trial functions do not satisfy the boundary condition. The Lagrange functional is

$$\mathfrak{L}(u) = J(u) + \int_0^1 \lambda_1^{(s)}(x) u(x, 0) dx + \int_0^1 \lambda_2^{(s)}(y) u(1, y) dy$$

$$-\int_{0}^{1} \lambda_{3}^{(s)}(x) u(x, 1) dx$$

$$-\int_{0}^{1} \lambda_{4}^{(s)}(y) u(0, y) dy,$$
(13)

where for convenience we split the Lagrange multiplier $\lambda^{(s)}$ into its restrictions $\lambda_i^{(s)}$ to each edge of the plate. This way we can represent the set of the boundary weight functions as the union of four sets selected separately for each edge; namely, $\lambda_i^{(s)}(z) = \sum_{j=1}^s \lambda_{i,j} \cos((j-1)\pi z)$, where $\lambda_{i,j}$ are the unknown coefficients. The number of internal equations here is again N^2 , and the number of the boundary equations is 4s, so the nondegeneracy condition is $N^2 > 4s$.

The exact solution to this problem can be expressed as a rapidly convergent series; we use its first ten terms to calculate the errors. We do not tabulate them here, because they are very large (up to 70%), and increasing the number of terms does not improve the approximation. At this point, this is not surprising since we are dealing with the same completeness issue as in the second example of Section 2. As we know, the system of cosines is incomplete on the interval, so the system of their products naturally is incomplete on the product of intervals that represents the plate. A more intuitive explanation is that products of cosines have vanishing normal derivatives $\partial u/\partial n$ on all edges of the plate, and all their linear combinations inherit this property. This does not matter for second-order equations because the normal derivatives are discontinuous on the relevant spaces, but it does matter for the higher-order equations like the biharmonic equation.

The choice of cosine products unwittingly enforces an additional boundary condition, $\partial u/\partial n=0$ on $\partial\Omega$. We therefore ended up solving a different variational problem. Together with u=0 on $\partial\Omega$ this describes, physically, a plate clamped on all sides (C-C-C-C) rather than a simply supported one. One can also check that in the weak formulation the boundary terms that multiply the variation of the solution's derivatives get removed because of the vanishing normal derivatives. Thus, we should be comparing our Ritz-Lagrange solutions to the answers for the C-C-C-C plate. Unfortunately, an analytic solution for a plate clamped on all sides is not known, but one can use the values obtained in [26, VI.44] numerically to make the comparison. The relative errors as percentages of the maximum deflection at the center of the plate are shown in Table 2.

To solve the original problem, we just need to complete the system of cosine products. By Lemma 7 it suffices to complete the cosines on the interval and take the products of the completed system. Namely, we take the products of x, x^2 , $\cos((i-1)\pi x)$ and y, y^2 , $\cos((i-1)\pi y)$ as the new trial functions and keep the rest of the above setting intact. The relative errors for the Ritz-Lagrange solutions with the completed system against the known series solution [2, 8.2.4] are shown in Table 3 and demonstrate the validity of the method.

As a final demonstration, we apply the Ritz-Lagrange method to a boundary eigenvalue problem for square plates. The eigenmodes describe standing vibrations of a plate, and

Table 2: Central (x = 0.5; y = 0.5) and boundary (x = 0; y = 0) error relative to maximum bending deflection of an C-C-C square plate of unit size.

| | Central error% | Boundary error% |
|---------------|----------------|-----------------|
| N = 4, s = 2 | 52.76 | -50.92 |
| N = 6, s = 3 | -3.3 | -1.22 |
| N = 8, s = 4 | 0.063 | -1.39 |
| N = 10, s = 5 | -0.12 | -0.09 |

Table 3: Central (x = 0.5; y = 0.5) and boundary (x = 0; y = 0) error relative to maximum bending deflection of an SS-SS-SS square plate of unit size.

| | Central error% | Boundary error% |
|---------------|----------------|-----------------|
| N = 5, s = 2 | 26.14 | -89.03 |
| N = 6, s = 3 | 2.81 | -3.7 |
| N = 8, s = 4 | 1.68 | -4.4 |
| N = 10, s = 5 | 0.56 | -1.1 |

their zeros (nodal curves) are known as Chladni figures [27, 5.1]. The problem has attracted a lot of attention from both analytic and numerical viewpoints; indeed Ritz himself applied his method to it in his original paper.

Boundary eigenvalue problems are somewhat beyond the scope of the theory in Section 3, which deals with linear constraints only. Under the Rayleigh-Ritz approach to solve for the eigenmodes one needs to impose an additional normalization constraint $(1/2) \int_{\Omega} u^2 dx = 1$ [23, 18.5], [24, VI.1.1], and [27, 5.2], which is quadratic. However, the general approach of Section 3 remains valid, and one can justify applying the Ritz-Lagrange method to problems with nonlinear constraints along the same lines.

Consider a uniformly loaded isotropic square plate of constant thickness with unit edge length, simply supported on all sides. The potential energy J of the plate is given by (12) without the distributed load term at the end. The boundary eigenvalue problem can be interpreted as finding extrema of J(u) subject to the boundary condition u = 0 on $\partial\Omega$, and the normalization constraint $(1/2) \int_{\Omega} u^2 dx = 1$.

Compared to (12) the Lagrange functional acquires an additional term $\mu(\int_{\Omega} u^2 dx - 1)$ and an additional equation, which amounts to the normalization constraint on the eigenmodes. When solving for trial solutions one can ignore this equation and use standard methods for finding eigenvectors instead. We keep the same choices for the trial functions and the boundary weights as before. Let a denote the vector of internal coefficients a_{ij} and let λ denote the vector of boundary coefficients $\lambda_{i,j}$. In terms of a and λ the Lagrange functional can be conveniently represented as $\mathfrak{L}(u^{(N)})$ = $(1/2)a^TKa - \mu((1/2)a^TMa - 1) + \lambda^TLa$, where K and M are matrices of size $N^2 \times N^2$ obtained by integrating the internal trial functions, see [2, 8.2.7], and L is a $4s \times N^2$ matrix obtained by integrating the boundary weights. Matrix L can be obtained by multiplying the boundary equations with the corresponding Lagrange multiplier functions and extracting the coefficients of a_{ij} and $\lambda_{i,j}$ after the integration. Note that

the boundary equations can be written as La = 0. Finally, differentiating the Lagrangian with respect to a_{ij} and $\lambda_{i,j}$, we are led to the following generalized eigenvalue problem:

$$\left(\left[\frac{K \mid L^{T}}{L \mid 0} \right] - \mu \left[\frac{M \mid 0}{0 \mid 0} \right] \right) \begin{pmatrix} a \\ \lambda \end{pmatrix} = 0. \tag{14}$$

For this eigenvalue problem to be solvable one needs L to have the maximal rank 4s, which is ensured by the nondegeneracy condition $N^2 \gg 4s$. The eigenvalues $\mu_i = \omega_i^2$ approximate the squares of the dimensionless natural frequencies of the plate's vibrations.

With N=10 and s=5 we obtain a set of approximate natural frequencies ω_i , first nine of which are shown below. Since the eigenmodes are known to be of the form $\sin(\pi mx)\sin(\pi ny)$ we change the single index notation to ω_{mn} and arrange the frequencies in a square pattern:

The exact values are taken from [2, 8.2.4]; the repeated frequencies correspond to multiple eigenvalues with the eigenmodes symmetric along different axes:

One can see that the estimated frequencies are slightly lower than the exact ones. This is in contrast with the application of the classical Ritz method, where the estimated frequencies are always higher. From a physical viewpoint, the latter happens because replacing an infinite system with a finite one is equivalent to imposing additional constraints, which tend to raise the stiffness of the system, and hence the frequencies. This assumes however that all the boundary constraints are enforced in both systems, that is, that the trial functions satisfy the essential boundary conditions.

In the Ritz-Lagrange method the trial functions do not satisfy the essential conditions, and even the trial solutions are forced to satisfy them only approximately. In other words, we are effectively relaxing the boundary constraints in addition to imposing additional ones through discretization. This relaxation lowers the frequencies (because a plate with fewer constraints is less stiff) and counteracts the effects of discretization. If we were able to impose the boundary conditions everywhere along the boundary the estimated frequencies would have been higher than the exact ones just as in the usual Ritz method.

From a mathematical viewpoint, this effect is also natural since the eigenvalues are the minima of a quadratic functional on subspaces of the original space [24, VI.1.1]. In the Ritz-Lagrange method we approximate them by using functions from a larger space (by relaxing the boundary conditions), thus lowering the minima that can be attained. In particular, one can see from the proof of Lemma 5 that the values of the

functional at the approximating minimizers are potentially smaller than at the true minimizer.

5. Conclusions

We developed a general extension of the Ritz method to systems of trial functions that do not satisfy the essential boundary conditions and proved its convergence. Our approach is based on treating the essential conditions as variational constraints and removing them using the Lagrange multipliers closely following the intuition behind the naive use of Lagrange multipliers. It is also more elementary than the Babuška-Brezzi saddle point formulation and leads to more transparent convergence conditions in terms of completeness of trial functions and boundary weights in appropriate spaces. We list some general observations on the workings of the method.

- (i) The variational functional has to be well-behaved not only on the energy space of the problem, but also on its extension that contains the trial functions. Sufficiently good behavior is a strong form of convexity, which in the case of quadratic functionals means that the boundary value problem is strongly elliptic.
- (ii) The systems of trial functions must be complete in the norm consistent with the functional, which is usually an extension of the energy norm of the problem to a larger space containing the trial functions. Although similar requirement applies to the classical Ritz method, it is much easier to encounter systems that appear complete but are not due to effects at the boundary.
- (iii) The Lagrange multipliers have to be treated as additional variables in the approximating systems. They can not be eliminated by substituting the trial solutions into the variational formulas for them in terms of the exact solution. These formulas are discontinuous in the relevant norms.
- (iv) In multidimensional problems the boundary conditions incorporate infinitely many constraints, and to obtain a finite-dimensional approximating system one has to select boundary weight functions in addition to the trial functions. The number of trial functions has to be significantly larger than the number of the boundary weights; otherwise the approximating system may be inconsistent or only have the trivial solution.
- (v) In multidimensional problems the approximating values of the functional may approach the exact value from below rather than above, in contrast to the classical Ritz method, because the minimization takes place on a larger space of functions not satisfying the boundary conditions exactly.
- (vi) The method can be applied to boundary eigenvalue problems interpreted (similar to the Rayleigh-Ritz approach) as minimization problems on subspaces of the original space with the additional normalization

constraint. Due to the presence of Lagrange multiplier variables the resulting finite-dimensional problem is a generalized eigenvalue problem $(A-\mu B)x=0$ instead of the ordinary one with B=I. In multidimensional vibrational problems the approximate eigenfrequencies obtained in this way may be lower than the exact ones (for the Ritz method they are always higher), due to relaxation of the boundary constraints.

As is well-known [2, 4], the Ritz method leads to the same approximating systems as the Galerkin method, but the latter can also be applied to nonoptimization problems. It is interesting if one can develop a "Galerkin-Lagrange method" without resorting to the saddle point formulation of Babuška-Brezzi. We believe that analogs of Theorem 6 can be proved for nonvariational equations with monotone operators using the approach of [16, VII.23.6] or [21, 26.2].

Appendix

Proofs

Proof of Lemma 1. Let $\{\phi_i\}$ be a complete system in \mathcal{U} . Since Γ has an *s*-dimensional image we can represent it as $\Gamma u = (\langle \Gamma_1, u \rangle, \dots, \langle \Gamma_s, u \rangle)^T$, where Γ_i are bounded linear functionals. Assume without loss of generality that they are linearly independent; otherwise some of them can be dropped without changing \mathcal{U} . Set $\mathcal{U}_k := \{u \in \mathcal{U} \mid \langle \Gamma_1, u \rangle, \dots, \langle \Gamma_k, u \rangle = 0\}$; we will construct a complete system in each \mathcal{U}_k by induction on k. Since $\mathcal{U} = \mathcal{U}_s$ the process concludes in *s* steps.

For k=1 we must produce a complete system of linear combinations in $\mathring{\mathcal{U}}_1 \coloneqq \{u \in \mathcal{U} \mid \langle \Gamma_1, u \rangle = 0\}$. Without loss of generality, $\langle \Gamma_1, \phi_1 \rangle \neq 0$ since $\{\phi_i\}$ is complete and Γ_1 can not vanish on all ϕ_i . We claim that $\widetilde{\phi}_i \coloneqq \phi_i - (\langle \Gamma_1, \phi_i \rangle / \langle \Gamma_1, \phi_1 \rangle) \phi_1 \in \mathring{\mathcal{U}}_1$ form the desired system. Let $u \in \mathring{\mathcal{U}}_1 \subset \mathcal{U}$ and a_i be the coefficients such that $\|u - \sum_{i=1}^N a_i \phi_i\| \leq \varepsilon$ for a given $\varepsilon > 0$. By definition of $\widetilde{\phi}_i$,

$$\sum_{i=1}^{N} a_i \widetilde{\phi}_i = \sum_{i=1}^{N} a_i \phi_i - \frac{\left\langle \Gamma_1, \sum_{i=1}^{N} a_i \phi_i \right\rangle}{\left\langle \Gamma_1, \phi_1 \right\rangle} \phi_1. \tag{A.1}$$

To estimate the second term we find,

$$\left| \left\langle \Gamma_{1}, \sum_{i=1}^{N} a_{i} \phi_{i} \right\rangle \right| = \left| \left\langle \Gamma_{1}, \sum_{i=1}^{N} a_{i} \phi_{i} - u \right\rangle + \left\langle \Gamma_{1}, u \right\rangle \right|$$

$$\leq \left\| \Gamma_{1} \right\| \left\| u - \sum_{i=1}^{N} a_{i} \phi_{i} \right\| \leq \left\| \Gamma_{1} \right\| \varepsilon.$$
(A.2)

Therefore, $\|u - \sum_{i=1}^{N} a_i \widetilde{\phi}_i\| \le (1 + \|\Gamma_1\| \|\phi_1\| / |\langle \Gamma_1, \phi_1 \rangle|) \varepsilon$, and since u, ε are arbitrary completeness of $\widetilde{\phi}_i$ follows.

Let $\{\tilde{\phi}_i\}$ be a complete system in $\mathring{\mathcal{U}}_k$ from the preceding step. Linear independence of Γ_j guarantees that Γ_{k+1} does not vanish on some $\tilde{\phi}_i$, which we may as well take to be $\tilde{\phi}_1$. Apply the process above with Γ_1 replaced by Γ_{k+1} and $\mathring{\mathcal{U}}_1$ replaced by

 $\mathring{\mathcal{U}}_{k+1}$ to obtain $\widetilde{\widetilde{\phi}}_i$. Then $\widetilde{\widetilde{\phi}}_i$ are linear combinations of $\widetilde{\phi}_i$ (and hence of the original ϕ_i), belong to $\mathring{\mathcal{U}}_{k+1}$, and are complete in it by the same argument. This concludes the induction step.

Proof of Theorem 2. A standard argument from convex analysis shows that if $J(u) \xrightarrow[\|u\| \to \infty]{} \infty$ on $\mathring{\mathscr{U}}$, then J has minimizers on $\mathring{\mathscr{U}}$, $\mathring{\mathscr{U}}^{(N)}$, and there is a weakly convergent subsequence $u^{(N_k)} \xrightarrow[k \to \infty]{} u^{(\infty)}$ [15, 6.2] and [16, III.10.3]; moreover $u^{(\infty)} \in \mathring{\mathscr{U}}$ since $0 = \Gamma(u^{(N_k)}) \xrightarrow[k \to \infty]{} \Gamma(u^{(\infty)})$. For large enough N there is a $u \in \mathring{\mathscr{U}}^{(N)}$ arbitrarily close to a minimizer \overline{u} of J on $\mathring{\mathscr{U}}$ by Lemma 1. Since J is continuous J(u) is arbitrarily close to the minimal value J_{\min} . But $J(u^{(N_k)})$ can not exceed J(u) for $N_k \geq N$ since $u^{(N_k)}$ is a minimizer on $\mathring{\mathscr{U}}^{(N_k)}$, so $J_{\min} \leq J(u^{(N_k)}) \leq J(u) = J_{\min} + \varepsilon$. Convex functionals are weakly lower semicontinuous [16, III.8.4], so after passing to limit we have that $J(u^{(\infty)}) \leq \lim_{k \to \infty} J(u^{(N_k)}) = J_{\min}$; that is, $u^{(\infty)}$ is a minimizer of J on $\mathring{\mathscr{U}}$ and $J(u^{(N_k)}) \xrightarrow[k \to \infty]{} J_{\min}$.

If we assume additionally that J is strictly convex on \mathcal{U} , then \overline{u} is unique and the entire sequence $u^{(N)}$ (which is now also uniquely defined) converges to it at least weakly [15, 6.2A].

The next two proofs use equivalent norms (inner products) on $W_2^1([0,\pi])$ and $W_2^2([0,\pi])$, respectively. Two norms are equivalent if they define the same notion of convergence; for equivalent norms on Sobolev spaces see [20, I.8] and especially [22, 1.9].

Proof of Lemma 3. The following inner product is equivalent to the usual one on $W_2^1([0,\pi])$: $\langle u,v\rangle_0:=u(0)v(0)+\int_0^\pi u_x v_x dx$. To prove completeness it suffices to show that any function w orthogonal to all cosines must be 0. For such w we have $\langle w,1\rangle_0=w(0)=0$ and hence $\langle w,\cos nx\rangle_0=\int_0^\pi w_x\cdot(-n\sin nx)dx=0$ for $n\geq 1$. Thus, w_x is L_2 orthogonal to $\sin nx$ for all $n\geq 1$. Since the latter form an orthogonal basis in $L_2([0,\pi])$ we must have $w_x=0$ a.e. But then by the Fundamental Theorem of Calculus $w(x)=w(0)+\int_0^x w_t dt=0$ a.e. establishing completeness. Being an orthogonal basis in $L_2([0,\pi])$ cosines must be minimal there, and therefore in any space with a stronger norm, which includes $W_2^1([0,\pi])$.

Proof of Lemma 4. An equivalent inner product on $W_2^2([0,\pi])$ is $\langle u,v\rangle_0 := u(0)v(0) + u_x(0)v_x(0) + \int_0^\pi u_{xx}v_{xx}dx$. Consider w orthogonal to all cosines; then we have $\langle w,1\rangle_0 = w(0) = 0$ and $\langle w,\cos nx\rangle_0 = \int_0^\pi w_{xx} \cdot (-n^2\cos nx)dx = 0$ for $n \ge 1$ because all sines vanish at 0. In particular, w_{xx} is L_2 orthogonal to $\cos nx$ for all $n \ge 1$. But orthogonal complement of the latter in L_2 consists of constants, so $w_{xx} = \text{const}$ and $w(x) = ax^2 + bx + c$. Since w(0) = 0 free term is 0 and w is a linear

combination of x and x^2 . Thus, orthogonal complement to cosines is spanned by x and x^2 proving completeness.

For minimality notice that by direct calculation $\langle x, \cos nx \rangle_0 = \langle x^2, \cos nx \rangle_0 = \langle x, x^2 \rangle_0 = 0$; that is, x and x^2 are orthogonal to all cosines and to each other. This means that neither one of them can be deleted without losing completeness. It also means that if a cosine can be approximated in W_2 by other cosines combined with x and x^2 , then it can already be approximated by other cosines alone. But the latter can not be done with arbitrary precision even in L_2 , let alone in W_2 .

Proof of Lemma 5. Since $\mathring{\mathcal{U}} \subset \cdots \subset \mathring{\mathcal{U}}_2 \subset \mathring{\mathcal{U}}_1$ and the minimum on a larger space can not get bigger we have $J(\overline{u}) \geq \cdots \geq J(\overline{u}_2) \geq J(\overline{u}_1)$. Thus, the numerical sequence $J(\overline{u}_s)$ is bounded. Moreover, $\overline{u}_s \in \mathring{\mathcal{U}}_{s_0}$ for $s \geq s_0$, so $\|\overline{u}_s\| \leq K < \infty$ for $s \geq s_0$ since J is weakly coercive on $\mathring{\mathcal{U}}_{s_0}$. The derivatives $J'(\overline{u}), J'(\overline{u}_s)$ vanish when paired with elements from subspaces where $\overline{u}, \overline{u}_s$, respectively, minimize J, so $\langle J'(\overline{u}), \overline{u} \rangle = \langle J'(\overline{u}_s), \overline{u} \rangle = \langle J'(\overline{u}_s), \overline{u}_s \rangle = 0$, and, by uniform monotonicity (8),

$$c\left(\left\|\overline{u}_{s}-\overline{u}\right\|\right) \leq \left\langle J'\left(\overline{u}_{s}\right)-J'\left(\overline{u}\right),\overline{u}_{s}-\overline{u}\right\rangle$$

$$=-\left\langle J'\left(\overline{u}\right),\overline{u}_{s}\right\rangle.$$
(A.3)

We prove below that the last expression converges to 0 when $s \to \infty$ implying that $\overline{u}_s \xrightarrow[s \to \infty]{} \overline{u}$ by norm since $c(t) \to 0$ implies $t \to 0$ by assumptions on c.

Since \overline{u} is a minimizer on \mathscr{U} the functional $J'(\overline{u})$ vanishes on any element from it. The subspace of functionals that vanish on the entire $\mathring{\mathscr{U}} = \{u \in \mathscr{U} \mid \Gamma u = 0\}$ is the closed linear span of $\{\Gamma^*\psi_j\}$ in \mathscr{U}^* . Indeed, if $\{\Gamma^*\psi_j\}$ did not span it there would exist, by the Hahn-Banach theorem, a u such that $\langle \Gamma^*\psi_j, u \rangle = \langle \psi_j, \Gamma u \rangle = 0$ for all j, while $\Gamma u \neq 0$, contradicting the completeness of $\{\psi_j\}$. Thus, for any $\varepsilon > 0$, there exists a linear combination $\xi = \sum_{j=1}^n a_j \Gamma^*\psi_j$ such that $\|J'(\overline{u}) - \xi\| \leq \varepsilon/K$. But then $\xi \in \mathring{\mathscr{U}}_n$, and for s > n we have $\langle \xi, \overline{u}_s \rangle = 0$, so

$$\left|\left\langle J'\left(\overline{u}\right),\overline{u}_{s}\right\rangle\right| = \left|\left\langle J'\left(\overline{u}\right) - \xi,\overline{u}_{s}\right\rangle\right| \leq \left\|J'\left(\overline{u}\right) - \xi\right\| \left\|\overline{u}_{s}\right\|$$

$$\leq \varepsilon.$$
(A.4)

Since ε is arbitrary $\langle J'(\overline{u}), \overline{u}_s \rangle \xrightarrow[s \to \infty]{} 0.$

Proof of Theorem 6. Let $\Gamma_s(u) := (\langle \psi_1, \Gamma u \rangle, \dots, \langle \psi_s, \Gamma u \rangle)^T$; then for $s \ge s_0$ the triple J, Γ_s , $\mathring{\mathcal{U}}_s$ satisfies conditions of Theorem 2. Moreover, uniform monotonicity of J' implies strict monotonicity, and hence strict convexity of J [16, III.5.3]. Therefore, by Theorem 2 there exist unique minimizers \overline{u}_s , $u_s^{(N)}$ on $\mathring{\mathcal{U}}_s$, $\mathring{\mathcal{U}}_s^{(N)}$, respectively, and $u_s^{(N)} \xrightarrow[N \to \infty]{w} \overline{u}_s$. The proof of strong convergence below relies on uniform monotonicity and more or less combines arguments from Theorem 23.3 in [16, VII.23.6] and Theorem 26.A(b) in [21, 26.2].

Since \overline{u}_s , $u_s^{(N)}$ are minimizers the derivatives $J'(u_s^{(N)})$, $J'(\overline{u}_s)$ vanish when paired with elements of $\mathring{\mathcal{U}}_s$, $\mathring{\mathcal{U}}_s$, respectively. Therefore, $\langle J'(\overline{u}_s), \overline{u}_s \rangle = \langle J'(\overline{u}_s), u_s^{(N)} \rangle = \langle J'(u_s^{(N)}), u_s^{(N)} \rangle = 0$, and by uniform monotonicity (8)

$$c\left(\left\|u_{s}^{(N)} - \overline{u}_{s}\right\|\right) \leq \left\langle J'\left(u_{s}^{(N)}\right) - J'\left(\overline{u}_{s}\right), u_{s}^{(N)} - \overline{u}_{s}\right\rangle$$

$$= -\left\langle J'\left(u_{s}^{(N)}\right), \overline{u}_{s}\right\rangle. \tag{A.5}$$

We prove below that $\langle J'(u_s^{(N)}), v \rangle \xrightarrow[N \to \infty]{} 0$ for any $v \in \mathring{\mathscr{U}}_s$, which implies $c(\|u_s^{(N)} - \overline{u}_s\|) \xrightarrow[N \to \infty]{} 0$, and hence $u_s^{(N)}$ converges to \overline{u}_s by norm.

First we show that the sequence of functionals $J'(u_s^{(N)})$ is uniformly bounded on $\mathring{\mathscr{U}}_s$. By monotonicity for any $h \in \mathring{\mathscr{U}}_s$ we have $\langle J'(u_s^{(N)}) - J'(h), u_s^{(N)} - h \rangle \geq 0$, so

$$\left\langle J'\left(u_{s}^{(N)}\right),h\right\rangle \leq \left\langle J'\left(u_{s}^{(N)}\right),u_{s}^{(N)}\right\rangle - \left\langle J'\left(h\right),u_{s}^{(N)}\right\rangle + \left\langle J'\left(h\right),h\right\rangle \tag{A.6}$$

$$= -\left\langle J'\left(h\right),u_{s}^{(N)}\right\rangle + \left\langle J'\left(h\right),h\right\rangle.$$

The right-hand side is clearly uniformly bounded for any $h \in \mathring{\mathcal{U}}_s$ since $u_s^{(N)}$ is weakly convergent, so $\|J'(u_s^{(N)})\| \leq M < \infty$ follows from the principle of uniform boundedness. Next for any $\varepsilon > 0$ choose N so large that $\|v - v_N\| \leq \varepsilon/M$ and $v_N \in \mathring{\mathcal{U}}_s^{(N)}$. This is possible due to completeness of $\{\phi_i\}$, and $\langle J'(u_s^{(N)}), v_N \rangle = 0$ since $u_s^{(N)}$ minimizes J on $\mathring{\mathcal{U}}_s^{(N)}$. Therefore,

$$\left| \left\langle J'\left(u_{s}^{(N)}\right), v \right\rangle \right| \leq \left| \left\langle J'\left(u_{s}^{(N)}\right), v - v_{N} \right\rangle \right|$$

$$\leq \left\| J'\left(u_{s}^{(N)}\right) \right\| \left\| v - v_{N} \right\| \leq \varepsilon,$$
(A.7)

and $\langle J'(u_s^{(N)}), v \rangle \xrightarrow[N \to \infty]{} 0$ as announced.

Convergence of \overline{u}_s to \overline{u} follows from Lemma 5, and convergence of $J(u_s^{(N)})$ and $J(\overline{u}_s)$ follows from the continuity of J.

Proof of Lemma 7. In the multi-index notation an equivalent norm on $W_p^k(\mathcal{D})$ is given by

$$||F||_{W_p^k} := \sum_{i=0}^k \sum_{|\alpha|=i} \left\| \frac{\partial^{|\alpha|} F}{\partial \xi^{\alpha}} \right\|_{L_p} = \sum_{|\alpha| \le k} \left\| \frac{\partial^{|\alpha|} F}{\partial \xi^{\alpha}} \right\|_{L_p}, \quad (A.8)$$

where the L_p norm is just $\|f\|_{L_p} \coloneqq (\int_{\mathscr{D}} |f|^p \, d\xi)^{1/p}$. If $f \in W_p^k(\Omega)$ and $\widetilde{f} \in W_p^k(\widetilde{\Omega})$ then it follows from the Fubini theorem that $\|f\widetilde{f}\|_{L_p} = \|f\|_{L_p} \|\widetilde{f}\|_{L_p}$ since f and \widetilde{f} depend on different variables. Let x and \widetilde{x} denote the variables on Ω

and $\widetilde{\Omega}$, respectively, so that $\xi = (x, \widetilde{x})$ is the variable on $\Omega \times \widetilde{\Omega}$. Then we estimate

$$\begin{split} \left\| F \widetilde{F} \right\|_{W_p^k} &= \sum_{|\alpha| \le k} \left\| \frac{\partial^{|\alpha|} F \widetilde{F}}{\partial \xi^{\alpha}} \right\|_{L_p} \\ &= \sum_{|\beta| + |\gamma| \le k} \left\| \frac{\partial^{|\beta|} F}{\partial x^{\beta}} \right\|_{L_p} \left\| \frac{\partial^{|\gamma|} \widetilde{F}}{\partial \widetilde{x}^{\gamma}} \right\|_{L_p} \\ &\le \sum_{|\beta| \le k, |\gamma| \le k} \left\| \frac{\partial^{|\beta|} F}{\partial x^{\beta}} \right\|_{L_p} \left\| \frac{\partial^{|\gamma|} \widetilde{F}}{\partial \widetilde{x}^{\gamma}} \right\|_{L_p} \\ &= \sum_{|\beta| \le k} \left\| \frac{\partial^{|\beta|} F}{\partial x^{\beta}} \right\|_{L_p} \sum_{|\gamma| \le k} \left\| \frac{\partial^{|\gamma|} \widetilde{F}}{\partial \widetilde{x}^{\gamma}} \right\|_{L_p} \\ &= \| F \|_{W_p^k} \left\| \widetilde{F} \right\|_{W^k}. \end{split} \tag{A.9}$$

Since ϕ_i are complete any monomial x^β can be approximated to any precision $\varepsilon>0$ in W_p^k by their linear combination $\phi=\sum_i a_i\phi_i$, and analogously \widetilde{x}^γ can be approximated by a linear combination $\widetilde{\phi}=\sum_j \widetilde{a}_j\widetilde{\phi}_j$. But $\phi\widetilde{\phi}=\sum_{i,j} a_i\widetilde{a}_j\phi_i\widetilde{\phi}_j$ is a linear combination of $\phi_i\widetilde{\phi}_j$, while the difference between the products can be made arbitrarily small:

$$\begin{aligned} \left\| x^{\beta} \widetilde{x}^{\gamma} - \phi \widetilde{\phi} \right\|_{W_{p}^{k}} &= \left\| x^{\beta} \left(\widetilde{x}^{\gamma} - \widetilde{\phi} \right) + \left(x^{\beta} - \phi \right) \widetilde{\phi} \right\|_{W_{p}^{k}} \\ &\leq \left\| x^{\beta} \right\|_{W_{p}^{k}} \left\| \widetilde{x}^{\gamma} - \widetilde{\phi} \right\|_{W_{p}^{k}} \\ &+ \left\| x^{\beta} - \phi \right\|_{W_{p}^{k}} \left\| \widetilde{\phi} \right\|_{W_{p}^{k}} \\ &\leq \varepsilon \left(\left\| x^{\beta} \right\|_{W_{p}^{k}} + \left\| \widetilde{x}^{\gamma} \right\|_{W_{p}^{k}} + \varepsilon \right), \end{aligned} \tag{A.10}$$

where the first inequality follows from the above estimate. Hence any product of monomials, and therefore any polynomial, can be approximated in W_p^k by linear combinations of $\phi_i \widetilde{\phi}_j$. By the generalized Weierstrass theorem [24, II.4.3], polynomials are complete in $W_p^k(\Omega \times \widetilde{\Omega})$ and hence so is the system $\{\phi_i \widetilde{\phi}_j\}$.

Competing Interests

The authors declare that they have no competing interests.

References

- [1] S. G. Mikhlin, *Variational Methods in Mathematical Physics*, Pergamon Press, New York, NY, USA, 1964.
- [2] J. N. Reddy, Energy Principles and Variational Methods in Applied Mechanics, John Wiley & Sons, Hoboken, NJ, USA, 2002.
- [3] L. Collatz, The Numerical Treatment of Differential Equations, Springer, Berlin, Germany, 1960.
- [4] H. H. E. Leipholz, "On some developments in direct methods of the calculus of variations," *Applied Mechanics Reviews*, vol. 40, no. 10, pp. 1379–1392, 1987.

- [5] C. Canuto, M. Y. Hussaini, A. Quarteroni, and T. A. Zang, Spectral Methods in Fluid Dynamics, Springer Series in Computational Physics, Springer, Berlin, Germany, 1988.
- [6] E. Gourgoulhon, "Introduction to spectral methods," in *Proceedings of the 4th EU Network Meeting*, Palma de Mallorca, Spain, September 2002, http://www.lorene.obspm.fr/palma.pdf.
- [7] W. Dahmen and A. Kunoth, "Appending boundary conditions by Lagrange multipliers: analysis of the LBB condition," *Numerische Mathematik*, vol. 88, no. 1, pp. 9–42, 2001.
- [8] I. Babuška, U. Banerjee, and J. E. Osborn, "Survey of meshless and generalized finite element methods: a unified approach," *Acta Numerica*, vol. 12, no. 12, pp. 1–125, 2003.
- [9] B. Budiansky and P. C. Hu, "The Lagrangian multiplier method of finding upper and lower limits to critical stresses of clamped plates," Tech. Rep. R-848, National Advisory Committee for Aeronautics, Langley Memorial Aeronautical Lab, 1946.
- [10] L. Klein, "Vibrations of constrained plates by a Rayleigh-Ritz method using lagrange multipliers," The Quarterly Journal of Mechanics & Applied Mathematics, vol. 30, no. 1, pp. 51–70, 1977.
- [11] I. Babuška, "The finite element method with Lagrangian multipliers," *Numerische Mathematik*, vol. 20, no. 3, pp. 179–192, 1973.
- [12] F. Brezzi, "On the existence, uniqueness and approximation of saddle-point problems arising from Lagrangian multipliers," *Revue Française d'Automatique, Informatique et Recherche Opérationnelle Série Rouge. Analyse Numérique*, vol. 8, no. 2, pp. 129–151, 1974.
- [13] F. Brezzi and M. Fortin, Mixed and Hybrid Finite Element Methods, Springer, New York, NY, USA, 1991.
- [14] C. Mollet, "Stability of Petrov-Galerkin discretizations: application to the space-time weak formulation for parabolic evolution problems," *Computational Methods in Applied Mathematics*, vol. 14, no. 2, pp. 231–255, 2014.
- [15] P. Drábek and J. Milota, Methods of Nonlinear Analysis. Applications to Differential Equations, Birkhäuser, Basel, Switzerland, 2007.
- [16] M. M. Vainberg, Variational Method and Method of Monotone Operators in the Theory of Nonlinear Equations, John Wiley & Sons, New York, NY, USA, 1973.
- [17] E. Zeidler, Nonlinear Functional Analysis and Its Applications— III. Variational Methods and Optimization, Springer, New York, NY, USA, 1985.
- [18] J. Storch and G. Strang, "Paradox lost: natural boundary conditions in the Ritz-Galerkin method," *International Journal for Numerical Methods in Engineering*, vol. 26, no. 10, pp. 2255–2266, 1988.
- [19] R. Goldberg, *Methods of Real Analysis*, John Wiley & Sons, New York, NY, USA, 1976.
- [20] O. A. Ladyzhenskaya, The Boundary Value Problems of Mathematical Physics, vol. 49 of Applied Mathematical Sciences, Springer, New York, NY, USA, 1985.
- [21] E. Zeidler, Nonlinear Functional Analysis and Its Applications— II/B. Nonlinear Monotone Operators, Springer, New York, NY, USA, 1990.
- [22] S. L. Sobolev, Some Applications of Functional Analysis in Mathematical Physics, vol. 90 of Translations of Mathematical Monographs, American Mathematical Society, Providence, RI, USA, 1991
- [23] E. Zeidler, Nonlinear Functional Analysis and Its Applications— II/A. Linear Monotone Operators, Springer, New York, NY, USA, 1990.

- [24] R. Courant and D. Hilbert, *Methods of Mathematical Physics I*, Interscience Publishers, New York, NY, USA, 1953.
- [25] R. Andreev, "Stability of sparse space-time finite element discretizations of linear parabolic evolution equations," *IMA Journal of Numerical Analysis*, vol. 33, no. 1, pp. 242–260, 2013.
- [26] S. Timoshenko and S. Woinowsky-Krieger, *Theory of Plates and Shells*, McGraw-Hill, Singapore, 1970.
- [27] M. J. Gander and G. Wanner, "From Euler, Ritz, and Galerkin to modern computing," SIAM Review, vol. 54, no. 4, pp. 627–666, 2012.



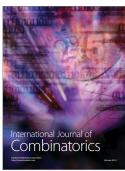








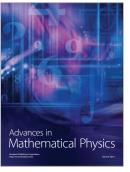






Submit your manuscripts at http://www.hindawi.com











Journal of Discrete Mathematics

