

Error Analysis of Nodal Meshless Methods

Robert Schaback

Abstract There are many application papers that solve elliptic boundary value problems by meshless methods, and they use various forms of generalized stiffness matrices that approximate derivatives of functions from values at scattered nodes $x_1, \dots, x_M \in \Omega \subset \mathbb{R}^d$. If u^* is the true solution in some Sobolev space S allowing enough smoothness for the problem in question, and if the calculated approximate values at the nodes are denoted by $\tilde{u}_1, \dots, \tilde{u}_M$, the canonical form of error bounds is

$$\max_{1 \leq j \leq M} |u^*(x_j) - \tilde{u}_j| \leq \epsilon \|u^*\|_S$$

where ϵ depends crucially on the problem and the discretization, but not on the solution. This contribution shows how to calculate such ϵ *numerically and explicitly*, for any sort of discretization of strong problems via nodal values, may the discretization use Moving Least Squares, unsymmetric or symmetric RBF collocation, or localized RBF or polynomial stencils. This allows users to compare different discretizations with respect to error bounds of the above form, without knowing exact solutions, and admitting all possible ways to set up generalized stiffness matrices. The error analysis is proven to be sharp under mild additional assumptions. As a byproduct, it allows to construct worst cases that push discretizations to their limits. All of this is illustrated by numerical examples.

1 Introduction

Following the seminal survey [5] by Ted Belytschko et.al. in 1996, meshless methods for PDE solving often work “*entirely in terms of values at nodes*”. This means that large linear systems are set up that have values $u(x_1), \dots, u(x_M)$ of an unknown

Robert Schaback
Institut für Numerische und Angewandte Mathematik, Univ. Göttingen, Lotzstraße 16-18, D-37083
Göttingen, Germany e-mail: schaback@math.uni-goettingen.de

function u as unknowns, while the equations model the underlying PDE problem in discretized way. Altogether, the discrete problems have the form

$$\sum_{j=1}^M a_{kj} u(x_j) \approx f_k, \quad 1 \leq k \leq N \quad (1)$$

with $N \geq M$, whatever the underlying PDE problem is, and the $N \times M$ matrix \mathbf{A} with entries a_{kj} can be called a *generalized stiffness matrix*.

Users solve the system somehow and then get values $\tilde{u}_1, \dots, \tilde{u}_M$ that satisfy

$$\sum_{j=1}^M a_{kj} \tilde{u}_j \approx f_k, \quad 1 \leq k \leq N,$$

but they should know how far these values are from the values $u^*(x_j)$ of the true solution of the PDE problem that is supposed to exist.

The main goal of this paper is to provide tools that allow users to assess the quality of their discretization, no matter how the problem was discretized or how the system was actually solved. The computer should tell the user whether the discretization is useful or not. It will turn out that this is possible, and at tolerable computational cost that is proportional to the complexity for setting up the system, not for solving it.

The only additional ingredient is a specification of the smoothness of the true solution u^* , and this is done in terms of a strong norm $\|\cdot\|_S$, e.g. a higher-order Sobolev norm or seminorm. The whole problem will then be implicitly scaled by $\|u^*\|_S$, and we assert an absolute bound of the form

$$\max_{1 \leq j \leq M} |u^*(x_j) - \tilde{u}_j| \leq \epsilon \|u^*\|_S$$

or a relative bound

$$\frac{\max_{1 \leq j \leq M} |u^*(x_j) - \tilde{u}_j|}{\|u^*\|_S} \leq \epsilon$$

with an entity ϵ that can be calculated. It will be a product of two values caring for *stability* and *consistency*, respectively, and these are calculated and analyzed separately.

Section 2 will set up the large range of PDE or, more generally, operator equation problems we are able to handle, and Section 3 provides the backbone of our error analysis. It must naturally contain some versions of *consistency* and *stability*, and we deal with these in Sections 5 and 7, with an interlude on polyharmonic kernels in Section 6. For given Sobolev smoothness order m , these provide stable, sparse, and error-optimal nodal approximations of differential operators. Numerical examples follow in Section 8, demonstrating how to work with the tools of this paper. It turns out that the evaluation of stability is easier than expected, while the evaluation of consistency often suffers from severe numerical cancellation that is to be overcome by future research, or that is avoided by using special scale-invariant approximations, e.g. via polyharmonic kernels along the lines of Section 6.

2 Problems and Their Discretizations

We have to connect the system (1) back to the original PDE problem, and we do this in an unconventional but useful way that we use successfully since [30] in 1999.

2.1 Analytic Problems

For example, consider a model boundary value problem of the form

$$\begin{aligned} Lu &= f \quad \text{in } \Omega \subset \mathbb{R}^d \\ Bu &= g \quad \text{in } \Gamma := \partial\Omega \end{aligned} \quad (2)$$

where f, g are given functions on Ω and Γ , respectively, and L, B are linear operators, defined and continuous on some normed linear space U in which the true solution u^* should lie. Looking closer, this is an infinite number of linear constraints

$$\begin{aligned} Lu(y) &= f(y) \quad \text{for all } y \in \Omega \subset \mathbb{R}^d \\ Bu(z) &= g(z) \quad \text{for all } z \in \Gamma := \partial\Omega \end{aligned}$$

and these can be generalized as infinitely many linear functionals acting on the function u , namely

$$\lambda(u) = f_\lambda \quad \text{for all } \lambda \in \Lambda \subset U^* \quad (3)$$

where the set Λ is contained in the topological dual U^* of U , in our example

$$\Lambda = \{\delta_y \circ L, y \in \Omega\} \cup \{\delta_z \circ B, z \in \Gamma\}. \quad (4)$$

Definition 1. An admissible problem in the sense of this paper consists in finding an u from some normed linear space U such that (3) holds for a fixed set $\Lambda \subset U^*$. Furthermore, solvability via $f_\lambda = \lambda(u^*)$ for all $\lambda \in \Lambda \subset U^*$ for some $u^* \in U$ is always assumed.

Clearly, this allows various classes of differential equations and boundary conditions, in weak or strong form. For examples, see [27]. Here, we just mention that the standard functionals for weak problems with $L = -\Delta$ are of the form

$$\lambda_v(u) := \int_{\Omega} (\nabla u)^T \nabla v \quad (5)$$

where v is an arbitrary test function from $W_0^1(\Omega)$.

2.2 Discretization

The connection of the problem (3) to the discrete linear system (1) usually starts with specifying a finite subset $\Lambda_N = \{\lambda_1, \dots, \lambda_N\} \subset \Lambda$ of *test* functionals. But then it splits into two essentially different branches.

The *shape function* approach defines functions $u_j : \Omega \rightarrow \mathbb{R}$ with the Lagrange property $u_i(x_j) = \delta_{ij}$, $1 \leq i, j \leq M$ and defines the elements a_{kj} of the stiffness matrix as $a_{kj} := \lambda_k(u_j)$. This means that the application of the functionals λ_k on *trial functions*

$$u(x) = \sum_{j=1}^M u(x_j) u_j(x)$$

is exact, and the linear system (1) describes the exact action of the selected test functionals on the trial space. Typical instances of the shape function approach are standard applications of Moving Least Squares (MLS) trial functions [32, 2, 3]. Such applications were surveyed in [5] and incorporate many versions of the Meshless Local Petrov Galerkin (MLPG) technique [4]. Another popular shape function method is unsymmetric or symmetric kernel-based collocation, see [16, 9, 11, 12].

But one can omit shape functions completely, at the cost of sacrificing exactness. Then the selected functionals λ_k are each approximated by linear combinations of the functionals $\delta_{x_1}, \dots, \delta_{x_M}$ by requiring

$$\lambda_k(u) \approx \sum_{j=1}^M a_{kj} \delta_{x_j} u = \sum_{j=1}^M a_{kj} u(x_j), \quad 1 \leq k \leq N, \quad \text{for all } u \in U. \quad (6)$$

This approach can be called *direct* discretization, because it bypasses shape functions. It is the standard technique for *generalized finite differences* (FD) [22], and it comes up again in meshless methods at many places, starting with [24, 31] and called *RBF-FD* or *local RBF collocation* by various authors, e.g. [34, 10]. The generalized finite difference approximations may be calculated via radial kernels using local selections of nodes only [25, 36, 35], and there are papers on how to calculate such approximations, e.g. [7, 18]. Bypassing Moving Least Squares trial functions, direct methods in the context of Meshless Local Petrov Galerkin techniques are in [21, 20], connected to *diffuse derivatives* [24]. For a mixture of kernel-based and MLS techniques, see [17].

This contribution will work in both cases, with a certain preference for the direct approach. The paper [27] focuses on shape function methods instead. It proves that uniform stability can be achieved for all well-posed problems by choosing a suitable discretization, and then convergence can be inferred from standard convergence rates of approximations of derivatives of the true solution from derivatives of trial functions. The methods of [27] fail for direct methods, and this was the main reason to write this paper.

2.3 Nodal Trial Approximations

In addition to Definition 1 we now assume that U is a space of functions on some set Ω , and that point evaluation is continuous, i.e. $\delta_x \in U^*$ for all $x \in \Omega$. We fix a finite set X_M of M nodes x_1, \dots, x_M and denote the span of the functionals δ_{x_j} by D_M .

For each $\lambda \in \Lambda$ we consider a linear approximation $\tilde{\lambda}$ to λ from D_M , i.e.

$$\lambda(u) \approx \sum_{j=1}^M a_j(\lambda)u(x_j) =: \tilde{\lambda}(u) \quad (7)$$

Note that there is no trial space of functions, and no *shape functions* at all, just *nodal values* and approximations of functionals from nodal values. It should be clear how the functionals in (4) can be approximated as in (7) via values at nodes.

In the sense of the preceding section, this looks like a *direct* discretization, but it also covers the *shape function* approach, because it is allowed to take $a_j(\lambda) = \lambda(u_j)$ for shape functions u_j with the Lagrange property.

2.4 Testing

Given a nodal trial approximation, consider a finite subset Λ_N of functionals $\lambda_1, \dots, \lambda_N$ and pose the possibly overdetermined linear system

$$\lambda_k(u^*) = f_{\lambda_k} = \sum_{j=1}^M a_j(\lambda_k)u_j \quad (8)$$

for unknown *nodal values* u_1, \dots, u_M that may be interpreted as approximations to $u^*(x_1), \dots, u^*(x_M)$. We call Λ_N a *test selection* of functionals, and remark that we have obtained a system of the form (1).

For what follows, we write the linear system (8) in matrix form als

$$\mathbf{f} = \mathbf{A}\mathbf{u} \quad (9)$$

with

$$\begin{aligned} \mathbf{A} &= (a_j(\lambda_k))_{1 \leq k \leq N, 1 \leq j \leq M} \in \mathbb{R}^{N \times M} \\ \mathbf{f} &= (f_{\lambda_1}, \dots, f_{\lambda_N})^T \in \mathbb{R}^N \\ \mathbf{u} &= (u(x_1), \dots, u(x_M))^T \in \mathbb{R}^M. \end{aligned}$$

Likewise, we denote the vector of exact nodal values $u^*(x_j)$ by \mathbf{u}^* , and $\tilde{\mathbf{u}}$ will be the vector of nodal values \tilde{u}_j that is obtained by some numerical method that solves the system (8) approximately.

It is well-known [14] that square systems of certain meshless methods may be singular, but it is also known [27] that one can bypass that problem by *overtesting*, i.e. choosing N larger than M . This leads to overdetermined systems, but they can be

handled by standard methods like the MATLAB backslash in a satisfactory way. Here, we expect that users set up their $N \times M$ stiffness matrix \mathbf{A} by sufficiently thorough testing, i.e. by selecting many test functionals $\lambda_1, \dots, \lambda_N$ so that the matrix has rank $M \leq N$. Section 7 will show that users can expect good stability if they handle a well-posed problem with sufficient overtesting. Note further that for cases like the standard Dirichlet problem (2), the set Λ_N has to contain a reasonable mixture of functionals connected to the differential operator and functionals connected to boundary values. Since we focus on general worst-case error estimates here, insufficient overtesting and an unbalanced mixture of boundary and differential equation approximations will result in error bounds that either cannot be calculated due to rank loss or come out large. The computer should reveal whether a discretization is good or not.

3 Error Analysis

The goal of this paper is to derive useful bounds for $\|\mathbf{u}^* - \tilde{\mathbf{u}}\|_\infty$, but we do not care for an error analysis away from the nodes. Instead, we assume a *postprocessing step* that interpolates the elements of $\tilde{\mathbf{u}}$ to generate an approximation \tilde{u} to the solution u^* in the whole domain. Our analysis will accept any numerical solution $\tilde{\mathbf{u}}$ in terms of nodal values and provide an error bound with small additional computational effort.

3.1 Residuals

We start with evaluating the *residual* $\mathbf{r} := \mathbf{f} - \mathbf{A}\tilde{\mathbf{u}} \in \mathbb{R}^N$ no matter how the numerical solution $\tilde{\mathbf{u}}$ was obtained. This can be explicitly done except for roundoff errors, and needs no derivation of upper bounds. Since in general the final error at the nodes will be larger than the observed residuals, users should refine their discretization when they encounter residuals that are very much larger than the expected error in the solution.

3.2 Stability

In Section 2.4 we postulated that users calculate an $N \times M$ stiffness matrix \mathbf{A} that has no rank loss. Then the *stability constant*

$$C_S(\mathbf{A}) := \sup_{\mathbf{u} \neq 0} \frac{\|\mathbf{u}\|_p}{\|\mathbf{A}\mathbf{u}\|_q} \quad (10)$$

is finite for any choice of discrete norms $\|\cdot\|_p$ and $\|\cdot\|_q$ on \mathbb{R}^M and \mathbb{R}^N , respectively, with $1 \leq p, q \leq \infty$ being fixed here, and dropped from the notation. In principle, this

constant can be explicitly calculated for standard norms, but we refer to Section 7 on how it is treated in theory and practice. We shall mainly focus on well-posed cases where $C_S(\mathbf{A})$ can be expected to be reasonably bounded, while norms of \mathbf{A} get very large. This implies that the ratios $\|\mathbf{u}\|_p/\|\mathbf{A}\mathbf{u}\|_q$ can vary in a wide range limited by

$$\|\mathbf{A}\|_{q,p}^{-1} \leq \frac{\|\mathbf{u}\|_p}{\|\mathbf{A}\mathbf{u}\|_q} \leq C_S(\mathbf{A}). \quad (11)$$

If we assume that we can deal with the stability constant $C_S(\mathbf{A})$, the second step of error analysis is

$$\begin{aligned} \|\mathbf{u}^* - \tilde{\mathbf{u}}\|_p &\leq C_S(\mathbf{A})\|\mathbf{A}(\mathbf{u}^* - \tilde{\mathbf{u}})\|_q \\ &\leq C_S(\mathbf{A})(\|\mathbf{A}\mathbf{u}^* - \mathbf{f}\|_q + \|\mathbf{f} - \mathbf{A}\tilde{\mathbf{u}}\|_q) \\ &\leq C_S(\mathbf{A})(\|\mathbf{A}\mathbf{u}^* - \mathbf{f}\|_q + \|\mathbf{r}\|_q) \end{aligned} \quad (12)$$

and we are left to handle the *consistency term* $\|\mathbf{A}\mathbf{u}^* - \mathbf{f}\|_q$ that still contains the unknown true solution \mathbf{u}^* . Note that \mathbf{f} is not necessarily in the range of \mathbf{A} , and we cannot expect to get zero residuals \mathbf{r} .

3.3 Consistency

For all approximations (7) we assume that there is a *consistency* error bound

$$|\lambda(u) - \tilde{\lambda}(u)| \leq c(\lambda)\|u\|_S \quad (13)$$

for all u in some *regularity* subspace U_S of U that carries a strong norm or seminorm $\|\cdot\|_S$. In case of a seminorm, we have to assume that the approximation $\tilde{\lambda}$ is an exact approximation to λ on the nullspace of the seminorm, but we shall use seminorms only in Section 6 below. If the solution u^* has plenty of smoothness, one may expect that $c(\lambda)\|u^*\|_S$ is small, provided that the discretization quality keeps up with the smoothness. In section 5, we shall consider cases where the $c(\lambda)$ can be calculated explicitly.

The bound (13) now specializes to

$$\|\mathbf{A}\mathbf{u}^* - \mathbf{f}\|_q \leq \|\mathbf{c}\|_q\|u^*\|_S$$

with the vector

$$\mathbf{c} = (c(\lambda_1), \dots, c(\lambda_N))^T \in \mathbb{R}^N,$$

and the error in (12) is bounded absolutely by

$$\|\mathbf{u}^* - \tilde{\mathbf{u}}\|_p \leq C_S(\mathbf{A}) \left(\|\mathbf{c}\|_q\|u^*\|_S + \|\mathbf{r}\|_q \right)$$

and relatively by

$$\frac{\|\mathbf{u}^* - \tilde{\mathbf{u}}\|_p}{\|u^*\|_S} \leq C_S(\mathbf{A}) \left(\|\mathbf{c}\|_q + \frac{\|\mathbf{r}\|_q}{\|u^*\|_S} \right). \quad (14)$$

This still contains the unknown solution u^* . But in kernel-based spaces, there are ways to get estimates of $\|u^*\|_S$ via interpolation. A strict but costly way is to interpolate the data vector \mathbf{f} by symmetric kernel collocation to get a function $u_{\mathbf{f}}^*$ with $\|u_{\mathbf{f}}^*\|_S \leq \|u^*\|_S$, and this norm can be plugged into (14). In single applications, users would prefer to take the values of $u_{\mathbf{f}}^*$ in the nodes as results, since they are known to be error-optimal [28]. But if discretizations with certain given matrices \mathbf{A} are to be evaluated or compared, this suggestion makes sense to get the right-hand side of (14) independent of u^* .

3.4 Residual Minimization

To handle the awkward final term in (14) without additional calculations, we impose a rather weak additional condition on the numerical procedure that produces $\tilde{\mathbf{u}}$ as an approximate solution to (9). In particular, we require

$$\|\mathbf{A}\tilde{\mathbf{u}} - \mathbf{f}\|_q \leq K(\mathbf{A})\|\mathbf{A}\mathbf{u}^* - \mathbf{f}\|_q, \quad (15)$$

which can be obtained with $K(\mathbf{A}) = 1$ if $\tilde{\mathbf{u}}$ is calculated via minimization of the residual $\|\mathbf{A}\mathbf{u} - \mathbf{f}\|_q$ over all $\mathbf{u} \in \mathbb{R}^M$, or with $K(\mathbf{A}) = 0$ if \mathbf{f} is in the range of \mathbf{A} . Anyway, we assume that users have a way to solve the system (9) approximately such that (15) holds with a known and moderate constant $K(\mathbf{A})$.

Then (15) implies

$$\begin{aligned} \|\mathbf{r}\|_q &= \|\mathbf{A}\tilde{\mathbf{u}} - \mathbf{f}\|_q \\ &\leq K(\mathbf{A})\|\mathbf{A}\mathbf{u}^* - \mathbf{f}\|_q \\ &\leq K(\mathbf{A})\|\mathbf{c}\|_q\|u^*\|_S \end{aligned}$$

and bounds $\|\mathbf{r}\|_q$ in terms of $\|u^*\|_S$.

3.5 Final Relative Error Bound

Theorem 1. *Under the above assumptions,*

$$\frac{\|\mathbf{u}^* - \tilde{\mathbf{u}}\|_p}{\|u^*\|_S} \leq (1 + K(\mathbf{A}))C_S(\mathbf{A})\|\mathbf{c}\|_q. \quad (16)$$

Proof. We can insert (15) directly into (12) to get

$$\begin{aligned} \|\mathbf{u}^* - \tilde{\mathbf{u}}\|_p &\leq C_S(\mathbf{A})(1 + K(\mathbf{A}))\|\mathbf{A}\mathbf{u}^* - \mathbf{f}\|_q \\ &\leq (1 + K(\mathbf{A}))C_S(\mathbf{A})\|\mathbf{c}\|_q\|u^*\|_S \end{aligned}$$

and finally (16), where now all elements of the right-hand side are accessible.

This is as far as one can go, not having any additional information on how u^* scales. The final form of (16) shows the classical elements of convergence analysis, since the right-hand side consists of a *stability* term $C_S(\mathbf{A})$ and a *consistency* term $\|\mathbf{c}\|_q$. The factor $1 + K(\mathbf{A})$ can be seen as a *computational accuracy* term.

Examples in Section 8 will show how these relative error bounds work in practice. Before that, the next sections will demonstrate theoretically why users can expect that the ingredients of the bound in (16) can be expected to be small. For this analysis, we shall assume that users know which regularity the true solution has, because we shall have to express everything in terms of $\|u^*\|_S$.

At this point, some remarks on error bounds should be made, because papers focusing on applications of meshless methods often contain one of the two standard crimes of error assessment.

The first is to take a problem with a known solution u^* that supplies the data, calculate nodal values $\tilde{\mathbf{u}}$ by some hopefully new method and then compare with \mathbf{u}^* to conclude that the method is good because $\|\mathbf{u}^* - \tilde{\mathbf{u}}\|$ is small. But the method may be intolerably unstable. If the input is changed very slightly, it may produce a seriously different numerical solution $\hat{\mathbf{u}}$ that reproduces the data as well as $\tilde{\mathbf{u}}$. The “quality” of the result $\tilde{\mathbf{u}}$ may be just lucky, it does not prove anything about the method used.

The second crime, usually committed when there is no explicit solution known, is to evaluate residuals $\mathbf{r} = \mathbf{A}\tilde{\mathbf{u}} - \mathbf{f}$ and to conclude that $\|\mathbf{u}^* - \tilde{\mathbf{u}}\|$ is small because residuals are small. This also ignores stability. There even are papers that claim convergence of methods by showing that residuals converge to zero when the discretization is refined. This reduces convergence rates of a PDE solver to rates of consistency, again ignoring stability problems that may counteract against good consistency. Section 8 will demonstrate this effect by examples.

This paper will avoid these crimes, but on the downside our error analysis is a worst-case theory that will necessarily overestimate errors of single cases.

3.6 Sharpness

In particular, if users take a specific problem (2) with data functions f and g and a known solution u^* , and if they evaluate the observed error and the bound (16), they will often see quite an overestimation of the error. This is due to the fact that they have a special case that is far away from being worst possible for the given PDE discretization, and this is comparable to a lottery win, as we shall prove now.

Theorem 2. *For all $K(\mathbf{A}) > 1$ there is some $u^* \in U_S$ and an admissible solution vector $\tilde{\mathbf{u}}$ satisfying (15) such that*

$$(K(\mathbf{A}) - 1)C_S(\mathbf{A})\|u^*\|_S\|\mathbf{c}\|_\infty \leq \|\mathbf{u}^* - \tilde{\mathbf{u}}\|_\infty \leq (K(\mathbf{A}) + 1)C_S(\mathbf{A})\|u^*\|_S\|\mathbf{c}\|_\infty \quad (17)$$

showing that the above worst-case error analysis cannot be improved much.

Proof. We first take the worst possible value vector \mathbf{u}_S for stability, satisfying

$$\|\mathbf{u}_S\|_\infty = C_S(\mathbf{A})\|\mathbf{A}\mathbf{u}_S\|_\infty$$

and normalize it to $\|\mathbf{u}_S\|_\infty = 1$. Then we consider the worst case of consistency, and we go into a kernel-based context.

Let the consistency vector \mathbf{c} attain its norm at some index j , $1 \leq j \leq N$, i.e. $\|\mathbf{c}\|_\infty = c(\lambda_j)$. Then there is a function $u_j \in U_S$ with

$$|\lambda_j(u_j) - \tilde{\lambda}_j(u_j)| = c(\lambda_j)\|u_j\|_S = c(\lambda_j)^2,$$

namely by taking the Riesz representer $u_j := (\lambda_j - \tilde{\lambda}_j)^x K(x, \cdot)$ of the error functional. The values of u_j at the nodes form a vector \mathbf{u}_j , and we take the data f as exact values of u_j , i.e. $f_k := \lambda_k(u_j)$, $1 \leq k \leq N$ to let u_j play the role of the true solution u^* , in particular $\mathbf{u}^* = \mathbf{u}_j$ and $\|u^*\|_S = \|u_j\|_S = c(\lambda_j) = \|\mathbf{c}\|_\infty$.

We then define $\tilde{\mathbf{u}} := \mathbf{u}^* + \alpha C_S(\mathbf{A})\mathbf{u}_S$ as a candidate for a numerical solution and check how well it satisfies the system and what its error bound is. We have

$$\begin{aligned} \|\mathbf{A}\tilde{\mathbf{u}} - \mathbf{f}\|_\infty &= \|\mathbf{A}(\mathbf{u}^* + \alpha C_S(\mathbf{A})\mathbf{u}_S) - \mathbf{f}\|_\infty \\ &\leq \|\mathbf{A}\mathbf{u}^* - \mathbf{f}\|_\infty + |\alpha|C_S(\mathbf{A})\|\mathbf{A}\mathbf{u}_S\|_\infty \\ &= |\alpha| + \|\mathbf{A}\mathbf{u}^* - \mathbf{f}\|_\infty \\ &= K(\mathbf{A})\|\mathbf{A}\mathbf{u}^* - \mathbf{f}\|_\infty \end{aligned}$$

if we choose

$$\alpha = (K(\mathbf{A}) - 1)\|\mathbf{A}\mathbf{u}^* - \mathbf{f}\|_\infty.$$

Thus $\tilde{\mathbf{u}}$ is a valid candidate for numerical solving. The actual error is

$$\begin{aligned} \|\mathbf{u}^* - \tilde{\mathbf{u}}\|_\infty &= (K(\mathbf{A}) - 1)\|\mathbf{A}\mathbf{u}^* - \mathbf{f}\|_\infty C_S(\mathbf{A}) \\ &= (K(\mathbf{A}) - 1)C_S(\mathbf{A}) \max_{1 \leq k \leq N} |\lambda_k(u_j) - \tilde{\lambda}_k(u_j)| \\ &\geq (K(\mathbf{A}) - 1)C_S(\mathbf{A})|\lambda_j(u_j) - \tilde{\lambda}_j(u_j)| \\ &= (K(\mathbf{A}) - 1)C_S(\mathbf{A})\|u_j\|_S\|\mathbf{c}\|_\infty \end{aligned} \quad (18)$$

proving the assertion.

We shall come back to this worst-case construction in the examples of Section 8.

4 Dirichlet Problems

The above error analysis simplifies for problems where Dirichlet values are given on boundary nodes, and where approximations of differential operators are only needed in interior points. Then we have N approximations of functionals that are based on M_I interior nodes and M_B boundary nodes, with $M = M_I + M_B$. We now use subscripts I and B to indicate vectors of values on interior and boundary nodes, respectively. The linear system now is

$$\mathbf{B}\mathbf{u}_I = \mathbf{f}_I - \mathbf{C}\mathbf{g}_B$$

while the previous section dealt with the full system

$$\mathbf{A} \begin{pmatrix} \mathbf{u}_I \\ \mathbf{u}_B \end{pmatrix} = \begin{pmatrix} \mathbf{f}_I \\ \mathbf{g}_B \end{pmatrix} \quad \text{with } \mathbf{A} = \begin{pmatrix} \mathbf{B} & \mathbf{C} \\ \mathbf{0} & \mathbf{I}_B \end{pmatrix}$$

that has trivial approximations on the boundary. Note that this splitting is standard practice in classical finite elements when nonzero Dirichlet boundary conditions are given. We now use the stability constant $C_S(\mathbf{B})$ for \mathbf{B} , not for \mathbf{A} , and examples will show that it often comes out much smaller than $C_S(\mathbf{A})$. The consistency bounds (13) stay the same, but they now take the form

$$\|\mathbf{B}\mathbf{u}_I^* + \mathbf{C}\mathbf{u}_B^* - \mathbf{f}_I\|_q = \|\mathbf{B}\mathbf{u}_I^* + \mathbf{C}\mathbf{g}_B - \mathbf{f}_I\|_q \leq \|\mathbf{c}_I\|_q \|u^*\|_S.$$

The numerical method should now guarantee

$$\|\mathbf{B}\tilde{\mathbf{u}}_I + \mathbf{C}\mathbf{g}_B - \mathbf{f}_I\|_q \leq K(\mathbf{B})\|\mathbf{B}\mathbf{u}_I^* + \mathbf{C}\mathbf{g}_B - \mathbf{f}_I\|_q$$

with a reasonable $K(\mathbf{B}) \geq 1$. Then the same error analysis applies, namely

$$\begin{aligned} \|\mathbf{u}_I^* - \tilde{\mathbf{u}}_I\|_p &\leq C_S(\mathbf{B})\|\mathbf{B}(\mathbf{u}_I^* - \tilde{\mathbf{u}}_I)\|_q \\ &\leq C_S(\mathbf{B})\|\mathbf{B}\mathbf{u}_I^* - \mathbf{C}\mathbf{g}_B - \mathbf{f}_I\|_q + C_S(\mathbf{B})\|\mathbf{B}\tilde{\mathbf{u}}_I - \mathbf{C}\mathbf{g}_B - \mathbf{f}_I\|_q \\ &\leq C_S(\mathbf{B})(1 + K(\mathbf{B}))\|\mathbf{B}\mathbf{u}_I^* - \mathbf{C}\mathbf{g}_B - \mathbf{f}_I\|_q \\ &\leq C_S(\mathbf{B})(1 + K(\mathbf{B}))\|\mathbf{c}_I\|_q \|u^*\|_S. \end{aligned}$$

5 Consistency Analysis

There are many ways to determine the *stiffness matrix elements* $a_j(\lambda_k)$ arising in (9) and (7), but they are either based on *trial/shape functions* or on *direct discretizations* as described in Section 2.2. We do not care here which technique is used. As a by-product, our method will allow to compare different approaches on a fair basis.

To make the constants $c(\lambda)$ in (13) numerically accessible, we assume that the norm $\|\cdot\|_S$ comes from a Hilbert subspace U_S of U that has a reproducing kernel

$$K : \Omega \times \Omega \rightarrow \mathbb{R}.$$

The squared norm of the error functional $\lambda - \tilde{\lambda}$ of the approximation $\tilde{\lambda}$ in (7) then is the value of the quadratic form

$$\begin{aligned}
Q^2(\lambda, \tilde{\lambda}) &:= \|\lambda - \tilde{\lambda}\|_{U_S^*}^2 \\
&= \lambda^x \lambda^y K(x, y) - 2 \sum_{j=1}^M a_j(\lambda) \lambda_j^x \lambda^y K(x, y) \\
&\quad + \sum_{j,k=1}^M a_j(\lambda) a_k(\lambda) \lambda_j^x \lambda_k^y K(x, y)
\end{aligned} \tag{19}$$

which can be explicitly evaluated, though there will be serious numerical cancellations because the result is small while the input is not. It provides the explicit error bound

$$|\lambda(u^*) - \tilde{\lambda}(u^*)|^2 \leq Q^2(\lambda, \tilde{\lambda}) \|u^*\|_S^2$$

such that we can work with

$$c(\lambda) = Q(\lambda, \tilde{\lambda}).$$

As mentioned already, the quadratic form (19) in its naïve form has an unstable evaluation due to serious cancellation. In [7], these problems were partly overcome by variable precision arithmetic, while the paper [18] provides a very nice stabilization technique, but unfortunately confined to approximations based on the Gaussian kernel. We hope to be able to deal with stabilization of the evaluation of the quadratic form in a forthcoming paper.

On the positive side, there are cases where these instabilities do not occur, namely for *polyharmonic kernels*. We shall come back to this in Section 6.

Of course, there are many *theoretical* results bounding the consistency error (13), e.g. [21, 7] in terms of $\|u^*\|_S$, with explicit convergence orders in terms of powers of *fill distances*

$$h := \sup_{y \in \Omega} \min_{x_j} \|y - x_j\|_2.$$

We call these orders *consistency orders* in what follows. Except for Section 6, we do not survey such results here, but users can be sure that a sufficiently fine fill distance and sufficient smoothness of the solution will always lead to a high consistency order. Since rates increase when more nodes are used, we target *p*-methods, not *h*-methods in the language of the finite element literature, and we assume sufficient regularity for this.

Minimizing the quadratic form (19) over the weights $a_j(\lambda)$ yields discretizations with *optimal* consistency with respect to the choice of the space U_S [7]. But their calculation may be unstable [18] and they usually lead to non-sparse matrices unless users restrict the used nodes for each single functional. If they are combined with a best possible choice of trial functions, namely the Riesz representers $v_j(x) = \lambda_j^y K(x, y)$ of the test functionals, the resulting linear system is symmetric and positive definite, provided that the functionals are linearly independent. This method is *symmetric collocation* [9, 11, 12], and it is an *optimal recovery* method in the space U_S [28]. It leads to non-sparse matrices and suffers from severe instability, but it is error-optimal. Here, we focus on non-optimal methods that allow sparsity.

Again, the instability of optimal approximations can be avoided using polyharmonic kernels, and the next section will describe how this works.

6 Approximations by Polyharmonic Kernels

Assume that we are working in a context where we know that the true solution u^* lies in Sobolev space $W_2^m(\Omega)$ for $\Omega \subset \mathbb{R}^d$, or, by Whitney extension also in $W_2^m(\mathbb{R}^d)$. Then the consistency error (13) of any given approximation should be evaluated in that space, and taking an optimal approximation in that space would yield a system with optimal consistency.

But since the evaluation and calculation of approximations in $W_2^m(\mathbb{R}^d)$ is rather unstable, a workaround is appropriate. Instead of the full norm in $W_2^m(\mathbb{R}^d)$ one takes the seminorm involving only the order m derivatives. This originates from early work of Duchon [8] and leads to Beppo-Levi spaces instead of Sobolev spaces (see e.g. [33]), but we take a summarizing shortcut here. Instead of the Whittle-Matérn kernel reproducing $W_2^m(\mathbb{R}^d)$, the radial *polyharmonic* kernel

$$H_{m,d}(r) := \begin{cases} (-1)^{\lceil m-d/2 \rceil} r^{2m-d}, & 2m-d \text{ odd} \\ (-1)^{1+m-d/2} r^{2m-d} \log r, & 2m-d \text{ even} \end{cases} \quad (20)$$

is taken, up to a scalar multiple

$$\begin{cases} \frac{\Gamma(m-d/2)}{2^{2m} \pi^{d/2} (m-1)!} & 2m-d \text{ odd} \\ \frac{1}{2^{2m-1} \pi^{d/2} (m-1)! (m-d/2)!} & 2m-d \text{ even} \end{cases} \quad (21)$$

that is used to match the seminorm in Sobolev space $W^m(\mathbb{R}^d)$. We allow m to be integer or half-integer. This kernel is *conditionally positive definite* of order $k = \lceil m - d/2 \rceil + 1$, and this has the consequence that approximations working in that space must be exact on polynomials of at least that order (= degree plus one). In some sense, this is the price to be paid for omitting the lower order derivatives in the Sobolev norm, but polynomial exactness will turn out to be a good feature, not a bug.

As an illustration for the connection between the polyharmonic kernel $H_{m,d}(r)$ and the Whittle-Matérn kernel $K_{m-d/2}(r)r^{m-d/2}$ reproducing $W_2^m(\mathbb{R}^d)$, we state the observation that (up to constants) the polyharmonic kernel arises as the first term in the expansion of the Whittle-Matérn kernel that is not an even power of r . For instance, up to higher-order terms,

$$K_3(r)r^3 = 16 - 2r^2 + \frac{1}{4}r^4 + \frac{1}{24}r^6 \log(r)$$

containing $H_{4,2}(r) = r^6 \log(r)$ up to a constant. This seems to hold in general for $K_n(r)r^n$ and $n = m - d/2$ for integer n and even dimension d . Similarly,

$$\frac{1}{\sqrt{2\pi}}K_{5/2}(r)r^{5/2} = 3 - \frac{1}{2}r^2 + \frac{1}{8}r^4 - \frac{1}{15}r^5$$

contains $H_{4,3}(r) = r^5$ up to a constant, and this generalizes to half-integer n with $n = m - d/2$. A rigid proof seems to be missing, but the upshot is that the polyharmonic kernel, if written with $r = \|x - y\|_2$, differs from the Whittle-Matérn kernel only by lower-order polynomials and higher-order terms, being simpler to evaluate.

If we have an arbitrary approximation (7) that is exact on polynomials of order k , we can insert its coefficients a_j into the usual quadratic form (19) using the polyharmonic kernel there, and evaluate the error. Clearly, the error is not smaller than the error of the optimal approximation using the polyharmonic kernel, and let us denote the coefficients of the latter by a_j^* .

We now consider *scaling*. Due to shift-invariance, we can assume that we have a homogeneous differential operator of order p that is to be evaluated at the origin, and we use scaled points hx_j for its nodal approximation. It then turns out [29] that the optimal coefficients $a_j^*(h)$ scale like $a_j^*(h) = h^{-p}a_j^*(1)$, and the quadratic form Q of (19) written in terms of coefficients as

$$\begin{aligned} Q^2(a) &= \lambda^x \lambda^y K(x, y) - 2 \sum_{j=1}^M a_j(\lambda) \lambda_j^x \lambda_j^y K(x, y) \\ &\quad + \sum_{j,k=1}^M a_j(\lambda) a_k(\lambda) \lambda_j^x \lambda_k^y K(x, y) \end{aligned}$$

scales *exactly* like

$$Q(a^*(h)) = h^{2m-d-2p} Q(a^*(1)),$$

proving that *there is no approximation of better order* in that space, no matter how users calculate their approximation. Note that strong methods (i.e. collocation) for second-order PDE problems (2) using functionals (4) have $p = 2$ while the weak functionals of (5) have $p = 1$. This is a fundamental difference between weak and strong formulations, but note that it is easy to have methods of arbitrarily high consistency order.

In practice, any set of given and centralized nodes x_j can be blown up to points Hx_j of average pairwise distance 1. Then the error and the weights can be calculated for the blown-up situation, and then the scaling laws for the coefficients and the error are applied using $h = 1/H$. This works for all scalings, without serious instabilities.

Now that we know an optimal approximation with a simple and stable scaling, why bother with other approximations? They will not have a smaller worst-case consistency error, and they will not always have the scaling property $a_j(h) = h^{-p}a_j(1)$, causing instabilities when evaluating the quadratic form. If they do have that scaling law, then

$$Q(a(h)) = h^{2m-d-2p} Q(a(1)) \geq h^{2m-d-2p} Q(a^*(1)) = Q(a^*(h))$$

can easily be proven, leading to stable calculation for an error that is not smaller than the optimal one. In contrast to standard results on the error of kernel-based

approximations, we have no restriction like $h \leq h_0$ here, since the scaling law is exact and holds for all h .

If the smoothness m for error evaluation is *fixed*, it will not pay off to use approximations with higher orders of polynomial exactness, or using kernels with higher smoothness. They cannot beat the optimal approximations for that smoothness class, and the error bounds of these are sharp. Special approximations can be better in a single case, but this paper deals with worst-case bounds, and then the optimal approximations are always superior.

The optimal approximations can be calculated for small numbers of nodes, leading to sparse stiffness matrices. One needs enough points to guarantee polynomial exactness of order $k = \lfloor m - d/2 \rfloor + 1$. The minimal number of points actually needed will depend on their geometric placement. The five-point star is an extremely symmetric example with exactness of order 4 in $d = 2$, but this order will normally need 15 points in general position because the dimension of the space of third-degree polynomials in \mathbb{R}^2 is 15.

The upshot of all of this is that, given a fixed smoothness m and a dimension d , polyharmonic stencils yield sparse optimal approximations that can be stably calculated and evaluated. Examples are in [29] and in Section 8 below. See [15] for an early work on stability of interpolation by polyharmonic kernels, and [1] for an example of an advanced application.

7 Stability Analysis

We now take a closer look at the stability constant $C_S(\mathbf{A})$ from (10). It can be rewritten as

$$C_S(\mathbf{A}) = \sup\{\|\mathbf{u}\|_p : \|\mathbf{A}\mathbf{u}\|_q \leq 1\} \quad (22)$$

and thus $2C_S(\mathbf{A})$ is the p -norm diameter of the convex set $\{\mathbf{u} \in \mathbb{R}^M : \|\mathbf{A}\mathbf{u}\|_q \leq 1\}$. In the case $p = q = \infty$ that will be particularly important below, this set is a polyhedron, and the constant $C_S(\mathbf{A})$ can be calculated via linear optimization. We omit details here, but note that the calculation tends to be computationally unstable and complicated. It is left to future research to provide a good estimation technique for the stability constant $C_S(\mathbf{A})$ like MATLAB's `condst` for estimating the L_1 condition number of a square matrix.

In case $p = q = 2$ we get

$$C_S(\mathbf{A})^{-1} = \min_{1 \leq j \leq M} \sigma_j$$

for the M positive singular values $\sigma_1, \dots, \sigma_M$ of A , and these are obtainable by *singular value decomposition*.

To simplify the computation, one might calculate the pseudoinverse \mathbf{A}^\dagger of \mathbf{A} and then take the standard (p, q) -norm of it, namely

$$\|\mathbf{A}^\dagger\|_{p,q} := \sup_{\mathbf{v} \neq 0} \frac{\|\mathbf{A}^\dagger \mathbf{v}\|_p}{\|\mathbf{v}\|_q}.$$

This overestimates $C_S(\mathbf{A})$ due to

$$\|\mathbf{A}^\dagger\|_{p,q} \geq \sup_{\mathbf{v}=\mathbf{A}\mathbf{u} \neq 0} \frac{\|\mathbf{A}^\dagger \mathbf{A}\mathbf{u}\|_p}{\|\mathbf{A}\mathbf{u}\|_q} = \sup_{\mathbf{u} \neq 0} \frac{\|\mathbf{u}\|_p}{\|\mathbf{A}\mathbf{u}\|_q} = C_S(\mathbf{A})$$

since $C_S(\mathbf{A})$ is the norm of the pseudoinverse not on all of \mathbb{R}^N , but restricted to the M -dimensional range of \mathbf{A} in \mathbb{R}^N . Here, we again used that \mathbf{A} has full rank, thus $\mathbf{A}^\dagger \mathbf{A} = I_{M \times M}$.

Calculating the pseudoinverse may be as expensive as the numerical solution of the system (8) itself, but if a user wants to have a close grip on the error, it is worth while. It assures stability of the numerical process, if not intolerably large, as we shall see. Again, we hope for future research to produce an efficient estimator.

A simple possibility, restricted to square systems, is to use the fact that MATLAB's `cond` estimates the 1-norm-condition number, which is the L_∞ condition number of the transpose. Thus

$$\tilde{C}_S(\mathbf{A}) := \frac{\text{cond}(\mathbf{A}')}{\|\mathbf{A}\|_\infty} \quad (23)$$

is an estimate of the L_∞ norm of \mathbf{A}^{-1} . This is computationally very cheap for sparse matrices and turns out to work fine on the examples in Section 8, but an extension to non-square matrices is missing.

We now switch to theory and want to show that users can expect $C_S(\mathbf{A})$ to be bounded above independent of the discretization details, if the underlying problem is well-posed. To this end, we use the approach of [27] in what follows.

Well-posed analytic problems of the form (3) allow a stable reconstruction of $u \in U$ from their full set of data $f_\lambda(u)$, $\lambda \in \Lambda$. This *analytic stability* can often be described as

$$\|u\|_{WP} \leq C_{WP} \sup_{\lambda \in \Lambda} |\lambda(u)| \text{ for all } u \in U, \quad (24)$$

where the *well-posedness norm* $\|\cdot\|_{WP}$ usually is weaker than the norm $\|\cdot\|_U$. For instance, elliptic second-order Dirichlet boundary value problems written in strong form satisfy

$$\|u\|_{\infty, \Omega} \leq \|u\|_{\infty, \partial\Omega} + C \|Lu\|_{\infty, \Omega} \text{ for all } u \in U := C^2(\Omega) \cap C(\bar{\Omega}), \quad (25)$$

see e.g. [6, (2.3), p. 14], and this is (24) for $\|\cdot\|_{WP} = \|\cdot\|_\infty$.

The results of [27] then show that for each trial space $U_M \subset U$ one can find a test set Λ_N such that (24) takes a discretized form

$$\|u\|_\infty \leq 2C_{WP} \sup_{\lambda_k \in \Lambda_N} |\lambda_k(u)| \text{ for all } u \in U_M,$$

and this implies

$$|u(x_j)| \leq 2C_{WP} \sup_{\lambda_k \in \Lambda_N} |\lambda_k(u)| \text{ for all } u \in U_M$$

for all nodal values. This proves a uniform stability property of the stiffness matrix with entries $\lambda_k(u_i)$. The functional approximations in [27] were of the form $a_j(\lambda) = \lambda(u_j)$, and then

$$\begin{aligned} \|\mathbf{u}\|_\infty &\leq 2C_{WP} \sup_{\lambda_k \in \Lambda_N} |\lambda_k(u)| \\ &= 2C_{WP} \sup_{\lambda_k \in \Lambda_N} |\lambda_k \left(\sum_{i=1}^M u(x_i) u_i \right)| \\ &= 2C_{WP} \sup_{\lambda_k \in \Lambda_N} \left| \sum_{i=1}^M u(x_i) \lambda_k(u_i) \right| \\ &= 2C_{WP} \|\mathbf{A}\mathbf{u}\|_\infty \end{aligned}$$

and thus

$$C_S(\mathbf{A}) \leq 2C_{WP}.$$

This is a prototype situation encouraging users to expect reasonably bounded norms of the pseudoinverse, provided that the norms are properly chosen.

However, the situation of [27] is much more special than here, because it is confined to the trial function approach. While we do not even specify trial spaces here, the paper [27] relies on the condition $a_j(\lambda) = \lambda(u_j)$ for a Lagrange basis of a trial space, i.e. exactness of the approximations on a chosen trial space. This is satisfied in nodal methods based on trial spaces, but not in direct nodal methods. In particular, it works for Kansa-type collocation and MLS-based nodal meshless methods, but not for localized kernel approximations and direct MLPG techniques in nodal form.

For general choices of $a_j(\lambda)$, the stability problem is a challenging research area that is not addressed here. Instead, users are asked to monitor the row-sum norm of the pseudoinverse numerically and apply error bounds like (16) for $p = q = \infty$. Note that the choice of discrete L_∞ norms is dictated by the well-posedness inequality (25). As pointed out above, chances are good to observe numerical stability for well-posed problems, provided that test functionals are chosen properly. We shall see this in the examples of Section 8. In case of square stiffness matrices, users can apply (23) to get a cheap and fairly accurate estimate of the stability constant.

For problems in weak form, the well-posedness norm usually is not $\|\cdot\|_{\infty, \Omega}$ but $\|\cdot\|_{L_2(\Omega)}$, and then we might get into problems using a nodal basis. In such cases, an L_2 -orthonormal basis would be needed for uniform stability, but we refrain from considering weak formulations here.

8 Examples

In all examples to follow, the nodal points are x_1, \dots, x_M in the domain $\Omega = [-1, +1]^2 \subset \mathbb{R}^2$, and parts of them are placed on the boundary. We consider the standard Dirichlet problem for the Laplacian throughout, and use testing points $y_1, \dots, y_n \in \Omega$ for the Laplacian and $z_1, \dots, z_k \in \partial\Omega$ for the Dirichlet boundary data in the sense of (4). Note that in our error bound (16) the right-hand sides of problems like (2) do not occur at all. This means that everything is only dependent on how the discretization works, it does not depend on any specific choice of f and g .

We omit detailed examples that show how the stability constant $C_S(\mathbf{A})$ decreases when increasing the number N of test functionals. An example is in [27], and (22) shows that stability must improve if rows are added to \mathbf{A} . Users are urged to make sure that their approximations (6), making up the rows of the stiffness matrix, have roughly the same consistency order, because adding equations will then improve stability without serious change of the consistency error.

We first take regular points on a 2D grid of sidelength h in $\Omega = [-1, +1]^2 \subset \mathbb{R}^2$ and interpret all points as nodes. On interior nodes, we approximate the Laplacian by the usual five-point star which is exact on polynomials up to degree 3 or order 4. On boundary nodes, we take the boundary values as given. This yields a square linear system. Since the coefficients of the five-point star blow up like $O(h^{-2})$ for $h \rightarrow 0$, the row-sum norm of \mathbf{A} and the condition must blow up like $O(h^{-2})$, which can easily be observed. The pseudoinverse does not blow up since the Laplacian part of \mathbf{A} just takes means and the boundary part is the identity. For the values of h we computed, its norm was bounded by roughly 1.3. This settles the stability issue from a practical point of view. Theorems on stability are not needed.

Consistency depends on the regularity space U_S chosen. We have a fixed classical discretization strategy via the five-point star, but we can evaluate the consistency error in different spaces. Table 1 shows the results for Sobolev space $W_2^4(\mathbb{R}^d)$. It clearly shows linear convergence, and its last column has the major part of the worst-case relative error bound (16). The estimate $\tilde{C}_S(\mathbf{A})$ from (23) agrees with $C_S(\mathbf{A})$ to all digits shown. Note that for all methods that need continuous point evaluations of the Laplacian in 2D, one cannot work with less smoothness, because the Sobolev inequality requires $W_2^m(\mathbb{R}^2)$ with $m > 2 + d/2 = 3$. The arguments in Section 6 show that the consistency order then is at most $m - d/2 - p = m - 3 = 1$, as observed. Table 2 shows the improvement if one uses the partial matrix \mathbf{B} of Section 4.

$M = N$	h	$C_S(\mathbf{A})$	$\ \mathbf{c}\ _\infty$	$C_S(\mathbf{A}) \ \mathbf{c}\ _\infty$
25	0.5000	1.281250	0.099045	0.126901
81	0.2500	1.291131	0.051766	0.066837
289	0.1250	1.293783	0.026303	0.034030
1089	0.0625	1.294459	0.013222	0.017116

Table 1 Results for five-point star on the unit square, for $W_2^4(\mathbb{R}^2)$ and the full matrix \mathbf{A}

$M_I = N_I$	h	$C_S(\mathbf{B})$	$\ \mathbf{c}_I\ _\infty$	$C_S(\mathbf{B}) \ \mathbf{c}_I\ _\infty$
9	0.5000	0.281250	0.099045	0.027856
49	0.2500	0.291131	0.051766	0.015071
225	0.1250	0.293783	0.026303	0.007727
961	0.0625	0.294459	0.013222	0.003893

Table 2 Results for five-point star on the unit square, for $W_2^4(\mathbb{R}^2)$ and the partial matrix \mathbf{B}

We now demonstrate the sharpness of our error bounds. We implemented the construction of Section 3.6 for $K(\mathbf{A}) = 2$ and the situation in the final row of Table 1. This means that, given \mathbf{A} , we picked values of f and g to realize worst-case stability and consistency, with known value vectors \mathbf{u}^* and $\tilde{\mathbf{u}}$. Figure 1 shows the values of \mathbf{u}_S and $\mathbf{u}_j = \mathbf{u}^*$ in the notation of the proof of Theorem 2, while Figure 2 displays $\tilde{\mathbf{u}}$. The inequality (17) is in this case

$$0.000226 = C_S(\mathbf{A})\|\mathbf{u}^*\|_S\|\mathbf{c}\|_\infty \leq \|\mathbf{u}^* - \tilde{\mathbf{u}}\|_\infty = 0.000226 \leq 3C_S(\mathbf{A})\|\mathbf{u}^*\|_S\|\mathbf{c}\|_\infty = 0.000679$$

and the admissibility inequality (15) is exactly satisfied with $K(\mathbf{A}) = 2$. Even though this example is worst-case, the residuals and the error $\|\mathbf{u}^* - \tilde{\mathbf{u}}\|_\infty$ are small compared to the last line of Table 1, and users might suspect that the table has a useless overestimation of the error. But the explanation is that the above bounds are absolute, not relative, while the norm of the true solution is $\|\mathbf{u}^*\|_S = \|\mathbf{c}\|_\infty = 0.0132$. The relative form of the above bound is

$$0.0171 = \frac{\|\mathbf{u}^* - \tilde{\mathbf{u}}\|_\infty}{\|\mathbf{u}^*\|_S} \leq 0.0513,$$

showing that the relative error bound 0.0171 in Table 1 is attained by a specific example. Thus our error estimation technique covers this situation well. The lower bound in the worst-case construction is attained because this example has equality in (18).

Note that our constructed case combines worst-case consistency with worst-case stability, but in practical situations these two worst cases will rarely happen at the same time. Figure 1 shows that the worst case for stability seems to be a discretization of a discontinuous function, and therefore it may be that practical situations are systematically far away from the worst case. This calls for a redefinition of the stability constant by restricting the range of \mathbf{A} in an appropriate way. The worst case for stability arises for vectors of nodal values that are close to the eigenvector of the smallest eigenvalue of \mathbf{A} , but the worst case for consistency might systematically have small inner products with eigenvectors for small eigenvalues.

If we take the polyharmonic kernel $H_{4,2}(r) = r^6 \log r$ (up to a constant), the five-point star is unique and therefore optimal, with consistency order 1, see Section 6. This means that for given smoothness order $m = 4$ and gridded nodes, the five-point star already has the optimal convergence order. Taking approximations of the Laplacian using larger subsets of nodes might be exact on higher-order polynomials, and will have smaller factors in front of the scaling law, but the consistency and convergence *order* will not be better, at the expense of losing sparsity.

To see how much consistency can be gained by using non-sparse optimal approximations by polyharmonic kernels, we worked at $h = 1$, approximating the error of the Laplacian at the origin by data in the integer nodes (m, n) with $-1 \leq m, n \leq K$ for increasing K . This models the case where the Laplacian is approximated in a near-corner point of the square. Smaller h can be handled by the scaling law. The consistency error in $W_2^4(\mathbb{R}^2)$ goes down from 0.07165 to 0.070035 when going from 25 to 225 neighbors (see Figure 3), while 0.08461 is the error of the five-point star

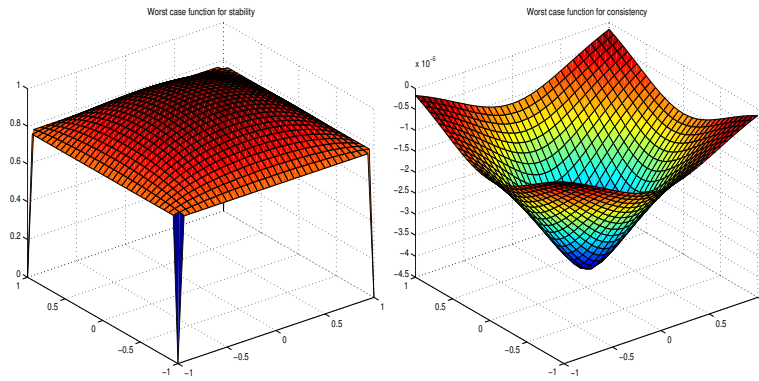


Fig. 1 Stability and consistency worst case

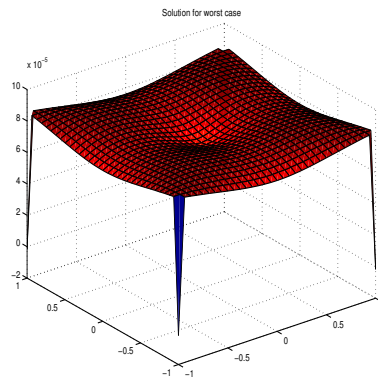


Fig. 2 Solution for joint worst case

at the origin. The gain is not worth the effort. The optimal stencils decay extremely quickly away from the origin. This is predicted by results of [19] concerning exponential decay of Lagrangians of polyharmonic kernels, as used successfully in [13] to derive local inverse estimates. See [23] for an early reference on polyharmonic near-Lagrange functions.

We now show how the technique of this paper can be used to compare very different discretizations, while a smoothness order m is fixed, in the sense that the true solution lies in Sobolev space $W_2^m(\Omega)$. Because we have p -methods in mind, we take $m = 6$ for the standard Dirichlet problem for the Laplacian in 2D and can expect an optimal consistency order $m - d/2 - 2 = 3$ for a strong discretization. Weak discretizations will be at least one order better, but we omit such examples. The required order of polynomial exactness when using the polyharmonic kernel is $1 + m - d/2 = 6$, which means that one should use at least 21 nodes for local approximations, if nodes are in general position, without symmetries. The bandwidth

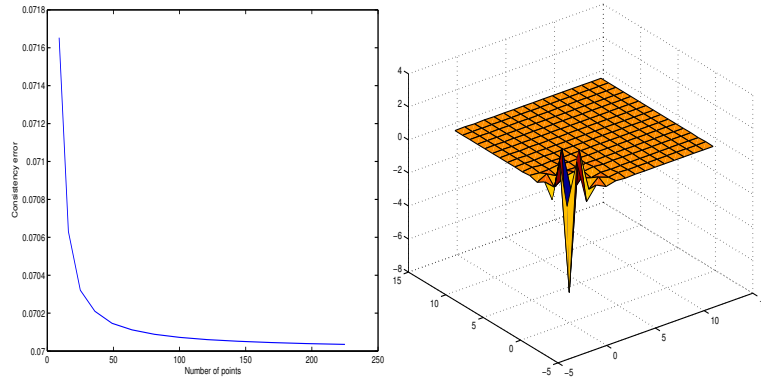


Fig. 3 Consistency error as a function of points offered, and stencil of optimal approximation for 225 nodes, as a function on the nodes

of the generalized stiffness matrix must therefore be at least 21. For convenience, we go to the unit square and a regular grid of meshwidth h first, to define the nodes. But then we add uniformly distributed noise of $\pm h/4$ to each interior node, keeping the boundary nodes. Then we approximate the Laplacian at each interior node locally by taking $n \geq 25$ nearest neighbor nodes, including boundary nodes, and set up the reduced generalized square stiffness matrix \mathbf{B} using the optimal polyharmonic approximation based on these neighboring nodes. On the boundary, we keep the given Dirichlet boundary values, following Section 4.

Table 3 shows results for local optimal approximations based on the polyharmonic kernel $H_{6,2}(r) = r^{10} \log r$ and $n = 30$ nearest neighbors. The stability constant was estimated via (23), for convenience and efficiency. One cannot expect to see an exact h^3 behavior in the penultimate column, since the nodes are randomly perturbed, but the overall behavior of the error is quite satisfactory. The computational complexity is roughly $O(Nn^3)$, and note that the linear system is not solved at all, because we used MATLAB's `condest`. Comparing with Table 4, it pays off to use a few more

$N = M$	$N_I = M_I$	h	$\tilde{C}_S(\mathbf{B})$	$\ \mathbf{e}_I\ _\infty$	$C_S(\mathbf{B})\ \mathbf{e}_I\ _\infty$
81	49	0.2500	2.3244	0.00075580	0.00175682
289	225	0.1250	0.3199	0.00005224	0.00001671
1089	961	0.0625	0.2964	0.00000872	0.00000259
4225	3969	0.0313	0.2961	0.00000147	0.00000044

Table 3 Optimal polyharmonic approximations using 30 neighbors

neighbors, and this also avoids instabilities. Users unaware of instabilities might think they can expect a similar behavior as in Table 3 when taking only 25 neighbors, but the third row of Table 4 should teach them otherwise. By resetting the random

number generator, all tables were made to work on the same total set of points, but the local approximations still yield rather different results.

The computationally cheapest way to calculate approximations with the required polynomial exactness of order 6 on 25 neighbors is to solve the linear 20×25 system describing polynomial exactness via the MATLAB backslash operator. It will return a solution based on 21 points only, i.e. with minimal bandwidth, but the overall behavior in Table 5 may not be worth the computational savings, if compared to the optimal approximations on 30 neighbors.

$N = M$	$N_I = M_I$	h	$\tilde{C}_S(\mathbf{B})$	$\ \mathbf{c}_I\ _\infty$	$C_S(\mathbf{B})\ \mathbf{c}_I\ _\infty$
81	49	0.2500	8.0180	0.00318328	0.02552351
289	225	0.1250	66.7176	0.00039055	0.02605641
1089	961	0.0625	417.8094	0.00003877	0.01620053
4225	3969	0.0313	75.5050	0.00000663	0.00050082

Table 4 Optimal polyharmonic approximations using 25 neighbors

$N = M$	$N_I = M_I$	h	$\tilde{C}_S(\mathbf{B})$	$\ \mathbf{c}_I\ _\infty$	$C_S(\mathbf{B})\ \mathbf{c}_I\ _\infty$
81	49	0.2500	9.0177	0.00354151	0.03193624
289	225	0.1250	25.6153	0.00058952	0.01510082
1089	961	0.0625	73.9273	0.00005482	0.00405249
4225	3969	0.0313	19.6458	0.00001186	0.00023305

Table 5 Backslash approximation on 25 neighbors

A more sophisticated kernel-based *greedy* technique [26, 29] uses between 21 and 30 points and works its way through the offered 30 neighbors to find a compromise between consistency error and support size. Table 6 shows the results, with an average of 23.55 neighbors actually used.

$N = M$	$N_I = M_I$	h	$\tilde{C}_S(\mathbf{B})$	$\ \mathbf{c}_I\ _\infty$	$C_S(\mathbf{B})\ \mathbf{c}_I\ _\infty$
81	49	0.2500	3.6188	0.00104016	0.00376411
289	225	0.1250	0.6128	0.00006821	0.00004180
1089	961	0.0625	0.3061	0.00000961	0.00000294
4225	3969	0.0313	0.2980	0.00000123	0.00000037

Table 6 Greedy polyharmonic approximations using at most 30 neighbors

For these examples, one can plot the consistency error as a function of the nodes, and there usually is a factor of 5 to 10 between the error in the interior and on the boundary. Therefore it should be better to let the node density increase towards the boundary, though this may lead to instabilities that may call for overtesting, i.e. to use $N \gg M$. For the same M and N as before, but with Chebyshev point distribution,

see Table 7. The additive noise on the interior points was 0.01, and we used the greedy method for up to 30 neighbors. This leads to a larger bandwidth near the corners, and to a consistency error that is now small at the boundary, see Figure 4. The average number of neighbors used was 23.3. Unfortunately, the scaling laws of stencils go down the drain here, together with the proven consistency order, but the results are still unexpectedly good.

$N = M$	$N_I = M_I$	h	$\tilde{C}_S(\mathbf{B})$	$\ \mathbf{c}_I\ _\infty$	$C_S(\mathbf{B})\ \mathbf{c}_I\ _\infty$
81	49	0.2500	111.1016	0.00433490	0.48161488
289	225	0.1250	0.4252	0.00006541	0.00002781
1089	961	0.0625	1.2133	0.00000677	0.00000821
4225	3969	0.0313	0.4353	0.00000120	0.00000052

Table 7 Greedy polyharmonic approximations using at most 30 neighbors, but in Chebyshev node arrangement

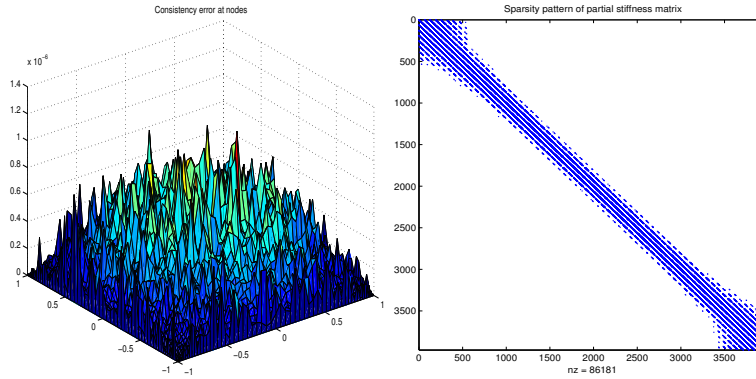


Fig. 4 Consistency plot and stiffness matrix \mathbf{B} for Chebyshev situation

For reasons of space and readability, we provide no examples for local approximations to weak functionals, and no comparisons with local approximations obtained via Moving Least Squares or the Direct Meshless Petrov Galerkin Method.

9 Conclusion and Outlook

The tables of the preceding section show that the numerical calculation of relative error bounds for PDE solving in spaces of fixed Sobolev smoothness can be done efficiently and with good results. This provides a general tool to evaluate discretizations in a worst-case scenario, without referring to single examples and complicated theorems.

Further examples should compare a large variety of competing techniques, the comparison being fair here as long as the smoothness m is fixed.

Users are strongly advised to use the cheap stability estimate (23) **anytime** to assess the stability of their discretization, if they have a square stiffness matrix. And, if they are not satisfied with the final accuracy, they should evaluate and plot the consistency error like in Figure 4 to see where the discretization should be refined. For all of this, polyharmonic kernels are an adequate tool.

It is left to future research to investigate and improve the stability estimation technique via (23), and, if the effort is worth while, to prove general theorems on sufficient criteria for stability. These will include assumptions on the placement of the trial nodes, as well as on the selection of sufficiently many and well-placed test functionals. In particular, stabilization by overtesting should work in general, but the examples in this paper show that overtesting may not be necessary at all. However, this paper serves as a practical workaround, as long as there are no theoretical cutting-edge results available.

Acknowledgements This work was strongly influenced by helpful discussions and e-mails with Oleg Davydov and Davoud Mirzaei.

References

1. T. Aboiyar, E.H. Georgoulis, and A. Iske, *Adaptive ADER methods using kernel-based polyharmonic spline WENO reconstruction*, SIAM Journal on Scientific Computing **32** (2010), 3251–3277.
2. M.G. Armentano, *Error estimates in Sobolev spaces for moving least square approximations*, SIAM J. Numer. Anal. **39**(1) (2001), 38–51.
3. M.G. Armentano and R.G. Durán, *Error estimates for moving least square approximations*, Appl. Numer. Math. **37** (2001), 397–416.
4. S. N. Atluri and T.-L. Zhu, *A new meshless local Petrov-Galerkin (MLPG) approach in Computational Mechanics*, Computational Mechanics **22** (1998), 117–127.
5. T. Belytschko, Y. Krongauz, D.J. Organ, M. Fleming, and P. Krysl, *Meshless methods: an overview and recent developments*, Computer Methods in Applied Mechanics and Engineering, special issue **139** (1996), 3–47.
6. D. Braess, *Finite elements. theory, fast solvers and applications in solid mechanics*, Cambridge University Press, 2001, Second edition.
7. O. Davydov and R. Schaback, *Error bounds for kernel-based numerical differentiation*, to appear in Numerische Mathematik, 2015.
8. J. Duchon, *Splines minimizing rotation-invariant semi-norms in Sobolev spaces*, Constructive Theory of Functions of Several Variables (W. Schempp and K. Zeller, eds.), Springer, Berlin–Heidelberg, 1979, pp. 85–100.
9. G. Fasshauer, *Solving partial differential equations by collocation with radial basis functions*, Surface Fitting and Multiresolution Methods (A. LeMéhauté, C. Rabut, and L.L. Schumaker, eds.), Vanderbilt University Press, Nashville, TN, 1997, pp. 131–138.
10. N. Flyer, E. Lehto, S. Blaise, G.B. Wright, and A. St.-Cyr, *A guide to RBF-generated finite differences for nonlinear transport: shallow water simulations on a sphere*, preprint, 2015.
11. C. Franke and R. Schaback, *Convergence order estimates of meshless collocation methods using radial basis functions*, Advances in Computational Mathematics **8** (1998), 381–399.

12. ———, *Solving partial differential equations by collocation using radial basis functions*, Appl. Math. Comp. **93** (1998), 73–82.
13. T. Hangelbroek, F.J. Narcowich, C. Rieger, and J.D. Ward, *An inverse theorem for compact Lipschitz regions using localized kernel bases*, arXiv preprint arXiv:1508.02952v2, 2015.
14. Y. C. Hon and R. Schaback, *On unsymmetric collocation by radial basis functions*, Appl. Math. Comput. **119** (2001), 177–186. MR MR1823674
15. A. Iske, *On the approximation order and numerical stability of local Lagrange interpolation by polyharmonic splines*, Modern Developments in Multivariate Approximation, Birkhäuser, Basel, 2003, pp. 153–165.
16. E. J. Kansa, *Application of Hardy’s multiquadric interpolation to hydrodynamics*, Proc. 1986 Simul. Conf., Vol. 4, 1986, pp. 111–117.
17. D. W. Kim and Y. Kim, *Point collocation methods using the fast moving least-square reproducing kernel approximation*, International Journal of Numerical Methods in Engineering **56** (2003), 1445–1464.
18. E. Larsson, E. Lehto, A. Heryodono, and B. Fornberg, *Stable computation of differentiation matrices and scattered node stencils based on Gaussian radial basis functions*, SIAM J. Sci. Comput. **35** (2013), A2096–A2119.
19. O.V. Matveev, *Spline interpolation of functions of several variables and bases in Sobolev spaces*, Trudy Mat. Inst. Steklov **198** (1992), 125–152.
20. D. Mirzaei and R. Schaback, *Direct Meshless Local Petrov-Galerkin (DMLPG) method: A generalized MLS approximation*, Applied Numerical Mathematics **68** (2013), 73–82.
21. D. Mirzaei, R. Schaback, and M. Dehghan, *On generalized moving least squares and diffuse derivatives*, IMA J. Numer. Anal. **32**, No. 3 (2012), 983–1000, doi: 10.1093/imanum/drr030.
22. A.R. Mitchell and D.F. Griffiths, *The finite difference method in partial differential equations*, John Wiley & Sons Ltd, 1980, pp. 233.
23. C. Rabut, *Elementary M -harmonic cardinal B -splines*, Numer. Algorithms **2** (1992), 39–62.
24. B. Nayroles, G. Touzot, and P. Villon, *Generalizing the finite element method: diffuse approximation and diffuse elements*, Computational Mechanics **10** (1992), 307–318.
25. B. Šarler, *From global to local radial basis function collocation method for transport phenomena*, Advances in meshfree techniques, Comput. Methods Appl. Sci., vol. 5, Springer, Dordrecht, 2007, pp. 257–282. MR 2433133 (2009i:65232)
26. R. Schaback, *Greedy sparse linear approximations of functionals from nodal data*, Numerical Algorithms **67** (2014), 531–547.
27. ———, *All well-posed problems have uniformly stable and convergent discretizations*, Numerische Mathematik **132** (2015), 243–269.
28. ———, *A computational tool for comparing all linear PDE solvers*, Advances of Computational Mathematics **41** (2015), 333–355.
29. ———, *Polyharmonic stencils*, Preprint, available from the author, 2015.
30. R. Schaback and H. Wendland, *Using compactly supported radial basis functions to solve partial differential equations*, Boundary Element Technology XIII (C. S. Chen, C. A. Brebbia, and D. W. Pepper, eds.), WitPress, Southampton, Boston, 1999, pp. 311–324.
31. A.I. Tolstykh, *On using radial basis functions in a “finite difference mode” with applications to elasticity problems*, Comput. Mech. **33** (2003), 68–79.
32. H. Wendland, *Local polynomial reproduction and moving least squares approximation*, IMA Journal of Numerical Analysis **21** (2001), 285–300.
33. ———, *Scattered data approximation*, Cambridge University Press, 2005.
34. G.B. Wright and B. Fornberg, *Scattered node compact finite difference-type formulas generated from radial basis functions*, J. Comput. Phys. **212** (2006), no. 1, 99–123. MR MR2183606 (2006j:65320)
35. G.M. Yao, B. Šarler, and C. S. Chen, *A comparison of three explicit local meshless methods using radial basis functions*, Eng. Anal. Bound. Elem. **35** (2011), no. 3, 600–609. MR 2753822 (2012b:65150)
36. G.M. Yao, Siraj ul Islam, and B. Šarler, *A comparative study of global and local meshless methods for diffusion-reaction equation*, CMES Comput. Model. Eng. Sci. **59** (2010), no. 2, 127–154. MR 2680809