

An adaptive stochastic Galerkin method

C.J. Gittelsohn

Research Report No. 2011-11
February 2011

Seminar für Angewandte Mathematik
Eidgenössische Technische Hochschule
CH-8092 Zürich
Switzerland

An Adaptive Stochastic Galerkin Method

Claude Jeffrey Gittelsohn*

February 28, 2011

Abstract

We derive an adaptive solver for random elliptic boundary value problems, using techniques from adaptive wavelet methods. Substituting wavelets by polynomials of the random parameters leads to a modular solver for the parameter dependence, which combines with any discretization on the spatial domain. We show optimality properties of this solver, and present numerical computations.

Introduction

Stochastic Galerkin methods have emerged in the past decade as an efficient solution procedure for boundary value problems depending on random data, see [DBO01, XK02, BTZ04, WK05, MK05, FST05, WK06, TS07, BS09, BAS10]. These methods approximate the random solution by a Galerkin projection onto a finite dimensional space of random fields. This requires the solution of a single coupled system of deterministic equations for the coefficients of the Galerkin projection with respect to a predefined set of basis functions on the parameter domain.

A major remaining obstacle is the construction of suitable spaces in which to compute approximate solutions. These should be adapted to the stochastic structure of the equation. Simple tensor product constructions are infeasible due to the high dimensionality of the parameter domain in case of input random fields with low regularity.

Parallel to but independently from the development of stochastic Galerkin methods, a new class of adaptive methods has emerged, which are set not in the continuous framework of a boundary value problem, but rather on the level of coefficients with respect to a hierarchic Riesz basis, such as a wavelet basis. Due to the norm equivalences constitutive of Riesz bases, errors and residuals in appropriate sequence spaces are equivalent to those in physically meaningful function spaces. This permits adaptive wavelet methods to be applied directly to a large class of equations, provided that a suitable Riesz basis is available.

For symmetric elliptic problems, the error of the Galerkin projection onto the span of a set of coefficients can be estimated using a sufficiently accurate approximation of the

*Research supported in part by the Swiss National Science Foundation grant No. 200021-120290/1.

residual of a previously computed approximate solution, see [CDD01, GHS07, DSS09]. This results in a sequence of finite-dimensional linear equations with successively larger sets of active coefficients.

We use techniques from these adaptive wavelet methods to derive an adaptive solver for random symmetric elliptic boundary value problems. In place of wavelets, we use an orthonormal polynomial basis on the parameter domain. The coefficients of the random solution with respect to this basis are deterministic functions on the spatial domain.

Adaptive wavelet methods extend to this vector setting, and lead to a modular solver which can be coupled with any discretization of or solver for the deterministic problem. We consider adaptive finite elements with a residual-based a posteriori error estimator.

We review random operator equations in Section 1. In particular, we derive the weak formulation of such equations, construct orthonormal polynomials on the parameter domain, and recast the weak formulation as a bi-infinite operator matrix equation for the coefficients of the random solution with respect to this polynomial basis. We refer to [Git11c] for further details.

A crucial ingredient in adaptive wavelet methods is the approximation of the residual. We study this for the setting of stochastic operator equations in Section 2. The resulting adaptive solver is presented in Section 3. We show convergence of the method, and provide a reliable error bound. Optimality properties are discussed in Section 4.

Finally, in Section 5, we apply the method to a simple elliptic equation. We discuss a suitable a posteriori finite element error estimator, and present numerical computations. These demonstrate the convergence of our solver and compare the adaptively constructed discretizations with the a priori adapted sparse tensor product construction from [BAS10]; we refer to [Git11b] for a comparison with other adaptive solvers. We discuss the empirical convergence behavior in the light of the theoretical approximation results in [CDS10b, CDS10a].

1 Stochastic Operator Equations

1.1 Pointwise Definition

Let $\mathbb{K} \in \{\mathbb{R}, \mathbb{C}\}$ and let V be a separable Hilbert space over \mathbb{K} . We denote by V^* the space of all continuous antilinear functionals on V . Furthermore, $\mathcal{L}(V, V^*)$ is the Banach space of bounded linear maps from V to V^* .

We consider operator equations depending on a parameter in $\Gamma := [-1, 1]^\infty$. Given

$$A: \Gamma \rightarrow \mathcal{L}(V, V^*) \quad \text{and} \quad f: \Gamma \rightarrow V^*, \quad (1.1)$$

we wish to determine

$$u: \Gamma \rightarrow V, \quad A(y)u(y) = f(y) \quad \forall y \in \Gamma. \quad (1.2)$$

Let $D \in \mathcal{L}(V, V^*)$ be the Riesz isomorphism, *i.e.* $\langle D \cdot, \cdot \rangle$ is the scalar product in V . We decompose A as

$$A(y) = D + R(y) \quad \forall y \in \Gamma \quad (1.3)$$

and assume that $R(y)$ is linear in $y \in \Gamma$,

$$R(y) = \sum_{m=1}^{\infty} y_m R_m \quad \forall y = (y_m)_{m=1}^{\infty} \in \Gamma, \quad (1.4)$$

as in *e.g.* [BAS10, BS09, CDS10b, CDS10a, TS07]. Here, each R_m is in $\mathcal{L}(V, V^*)$. We assume $(R_m)_m \in \ell^1(\mathbb{N}; \mathcal{L}(V, V^*))$, and there is a $\gamma \in [0, 1)$ such that $\|R(y)\|_{V \rightarrow V^*} \leq \gamma$ for all $y \in \Gamma$. By [Git11c, Proposition 1.2], this ensures existence and uniqueness of the solution of (1.1). For simplicity, we also assume that the sequence $(\|R_m\|_{V \rightarrow V^*})_{m=1}^{\infty}$ is nonincreasing.

1.2 Weak Formulation

Let π be a probability measure on the parameter domain Γ with Borel σ -algebra $\mathcal{B}(\Gamma)$. We assume that the map $\Gamma \ni y \mapsto A(y)v(y)$ is measurable for any measurable $v: \Gamma \rightarrow V$. Then

$$\mathcal{A}: L_{\pi}^2(\Gamma; V) \rightarrow L_{\pi}^2(\Gamma; V^*), \quad v \mapsto [y \mapsto A(y)v(y)], \quad (1.5)$$

is well-defined and continuous. We assume also that $f \in L_{\pi}^2(\Gamma; V^*)$.

The weak formulation of (1.2) is to find $u \in L_{\pi}^2(\Gamma; V)$ such that

$$\int_{\Gamma} \langle A(y)u(y), v(y) \rangle d\pi(y) = \int_{\Gamma} \langle f(y), v(y) \rangle d\pi(y) \quad \forall v \in L_{\pi}^2(\Gamma; V). \quad (1.6)$$

The left term in (1.6) is the duality pairing in $L_{\pi}^2(\Gamma; V)$ of $\mathcal{A}u$ with the test function v , and the right term is the duality pairing of f with v . We follow the convention that the duality pairing is linear in the first argument and antilinear in the second.

By [Git11c, Theorem 1.4], the solution u of (1.2) is in $L_{\pi}^2(\Gamma; V)$, and it is the unique solution of (1.6). In particular, the operator \mathcal{A} is boundedly invertible.

We define the multiplication operators

$$K_m: L_{\pi}^2(\Gamma) \rightarrow L_{\pi}^2(\Gamma), \quad v(y) \mapsto y_m v(y), \quad m \in \mathbb{N}. \quad (1.7)$$

Since y_m is real and $|y_m|$ is less than one, K_m is symmetric and has norm at most one.

By separability of V , the Lebesgue–Bochner space $L_{\pi}^2(\Gamma; V)$ is isometrically isomorphic to the Hilbert tensor product $L_{\pi}^2(\Gamma) \otimes V$, and similarly for V^* in place of V . Using these identifications, we expand \mathcal{A} as $\mathcal{A} = \mathcal{D} + \mathcal{R}$ with

$$\mathcal{D} := \text{id}_{L_{\pi}^2(\Gamma)} \otimes D \quad \text{and} \quad \mathcal{R} := \sum_{m=1}^{\infty} K_m \otimes R_m. \quad (1.8)$$

This sum converges in $\mathcal{L}(L_{\pi}^2(\Gamma; V), L_{\pi}^2(\Gamma; V^*))$ by the assumption $(R_m)_m \in \ell^1(\mathbb{N}; \mathcal{L}(V, V^*))$.

Lemma 1.1. $\|\mathcal{R}\|_{L_{\pi}^2(\Gamma; V) \rightarrow L_{\pi}^2(\Gamma; V^*)} \leq \gamma < 1$.

Proof. We note that, as in (1.5), $(\mathcal{R}v)(y) = R(y)v(y)$ for all $v \in L_{\pi}^2(\Gamma; V)$ and $y \in \Gamma$. Therefore, using the assumption $\|R(y)\|_{V \rightarrow V^*} \leq \gamma$,

$$\|\mathcal{R}v\|_{L_{\pi}^2(\Gamma; V^*)}^2 = \int_{\Gamma} \|R(y)v(y)\|_{V^*}^2 d\pi(y) \leq \int_{\Gamma} \|R(y)\|_{V \rightarrow V^*}^2 \|v(y)\|_V^2 d\pi(y) \leq \gamma^2 \|v\|_{L_{\pi}^2(\Gamma; V)}^2. \quad \square$$

1.3 Orthonormal Polynomial Basis

In order to construct an orthonormal polynomial basis of $L^2_\pi(\Gamma)$, we assume that π is a product measure. Let

$$\pi = \bigotimes_{m=1}^{\infty} \pi_m \quad (1.9)$$

for probability measures π_m on $([-1, 1], \mathcal{B}([-1, 1]))$; see *e.g.* [Bau02, Section 9] for a general construction of arbitrary products of probability measures. We assume that the support of π_m in $[-1, 1]$ has infinite cardinality.

For all $m \in \mathbb{N}$, let $(P_n^m)_{n=0}^{\infty}$ be an orthonormal polynomial basis of $L^2_{\pi_m}([-1, 1])$, with $\deg P_n^m = n$. Such a basis is given by the three term recursion $P_{-1}^m := 0$, $P_0^m := 1$ and

$$\beta_n^m P_n^m(\xi) := (\xi - \alpha_{n-1}^m) P_{n-1}^m(\xi) - \beta_{n-1}^m P_{n-2}^m(\xi), \quad n \in \mathbb{N}, \quad (1.10)$$

with

$$\alpha_n^m := \int_{-1}^1 \xi P_n^m(\xi)^2 d\pi_m(\xi) \quad \text{and} \quad \beta_n^m := \frac{c_{n-1}^m}{c_n^m}, \quad (1.11)$$

where c_n^m is the leading coefficient of P_n^m , $\beta_0^m := 1$, and $P_m n$ is chosen as normalized in $L^2_{\pi_m}([0, 1])$ with a positive leading coefficient. This basis is unique *e.g.* if c_n^m is chosen to be positive.

We define the set of finitely supported sequences in \mathbb{N}_0 as

$$\Lambda := \{v \in \mathbb{N}_0^{\mathbb{N}}; \#\text{supp } v < \infty\}, \quad (1.12)$$

where the support is defined by

$$\text{supp } v := \{m \in \mathbb{N}; v_m \neq 0\}, \quad v \in \mathbb{N}_0^{\mathbb{N}}. \quad (1.13)$$

Then countably infinite tensor product polynomials are given by

$$P := (P_v)_{v \in \Lambda}, \quad P_v := \bigotimes_{m=1}^{\infty} P_{v_m}^m, \quad v \in \Lambda. \quad (1.14)$$

Note that each of these functions depends on only finitely many dimensions,

$$P_v(y) = \prod_{m=1}^{\infty} P_{v_m}^m(y_m) = \prod_{m \in \text{supp } v} P_{v_m}^m(y_m), \quad v \in \Lambda, \quad (1.15)$$

since $P_0^m = 1$ for all $m \in \mathbb{N}$.

By *e.g.* [Git11c, Theorem 2.8], P is an orthonormal basis of $L^2_\pi(\Gamma)$. By Parseval's identity, this is equivalent to the statement that the map

$$T: \ell^2(\Lambda) \rightarrow L^2_\pi(\Gamma), \quad (c_v)_{v \in \Lambda} \mapsto \sum_{v \in \Lambda} c_v P_v, \quad (1.16)$$

is a unitary isomorphism. The inverse of T is

$$T^{-1} = T^*: L^2_\pi(\Gamma) \rightarrow \ell^2(\Lambda), \quad g \mapsto \left(\int_{\Gamma} g(y) \overline{P_v(y)} d\pi(y) \right)_{v \in \Lambda}. \quad (1.17)$$

1.4 Bi-Infinite Operator Matrix Equation

We use the isomorphism T from (1.16) to recast the weak stochastic operator equation (1.6) as an equivalent discrete operator equation. Since T is a unitary map from $\ell^2(\Lambda)$ to $L^2_\pi(\Gamma)$, the tensor product operator $T_V := T \otimes \text{id}_V$ is an isometric isomorphism from $\ell^2(\Lambda; V)$ to $L^2_\pi(\Gamma; V)$. By definition, $w \in L^2_\pi(\Gamma; V)$ and $\mathbf{w} = (w_\nu)_{\nu \in \Lambda} \in \ell^2(\Lambda; V)$ are related by $w = T_V \mathbf{w}$ if

$$w(y) = \sum_{\nu \in \Lambda} w_\nu P_\nu(y) \quad \text{or} \quad w_\nu = \int_\Gamma w(y) \overline{P_\nu(y)} d\pi(y) \quad \forall \nu \in \Lambda, \quad (1.18)$$

and either of these properties implies the other. The series in (1.18) converges unconditionally in $L^2_\pi(\Gamma; V)$, and the integral can be interpreted as a Bochner integral in V .

Let $A := T_V^* \mathcal{A} T_V$ and $f := T_V^* f$. Then $u = T_V u$ for $u \in \ell^2(\Lambda; V)$ with

$$A u = f \quad (1.19)$$

since $u \in L^2_\pi(\Gamma; V)$ satisfies $\mathcal{A} u = f$.

By definition, A is a boundedly invertible linear map from $\ell^2(\Lambda; V)$ to $\ell^2(\Lambda; V^*)$. It can be interpreted as a bi-infinite operator matrix

$$A = [A_{\nu\mu}]_{\nu, \mu \in \Lambda}, \quad A_{\nu\mu} : V \rightarrow V^*, \quad (1.20)$$

with entries

$$\begin{aligned} A_{\nu\nu} &= D + \sum_{m=1}^{\infty} \alpha_{\nu_m}^m R_m, \quad \nu \in \Lambda, \\ A_{\nu\mu} &= \beta_{\max(\nu_m, \mu_m)}^m R_m, \quad \nu, \mu \in \Lambda, \quad \nu - \mu = \pm \epsilon_m, \end{aligned} \quad (1.21)$$

and $A_{\nu\mu} = 0$ otherwise, where ϵ_m denotes the Kronecker sequence with $(\epsilon_m)_n = \delta_{mn}$. If π_m is a symmetric measure on $[-1, 1]$ for all $m \in \mathbb{N}$, then $\alpha_n^m = 0$ for all m and n , and thus $A_{\nu\nu} = D$. We refer to [Git11c, Git11a] for details.

Similarly, the operator $R := T_V^* \mathcal{R} T_V$ can be interpreted as a bi-infinite operator matrix $R = [R_{\nu\mu}]$ with $R_{\nu\nu} = A_{\nu\nu} - D$ and $R_{\nu\mu} = A_{\nu\mu}$ for $\nu \neq \mu$.

Let $K_m = T^* K_m T \in \mathcal{L}(\ell^2(\Lambda))$. Due to the three term recursion (1.10),

$$(\mathbf{K}_m \mathbf{c})_\mu = \beta_{\mu_m+1}^m c_{\mu+\epsilon_m} + \alpha_{\mu_m}^m c_\mu + \beta_{\mu_m}^m c_{\mu-\epsilon_m}, \quad \mu \in \Lambda, \quad (1.22)$$

for $\mathbf{c} = (c_\mu)_{\mu \in \Lambda} \in \ell^2(\Lambda)$, where $c_\mu := 0$ if $\mu_m < 0$ for any $m \in \mathbb{N}$. Furthermore, $\mathbf{K}_m^* = \mathbf{K}_m$ and $\|\mathbf{K}_m\|_{\ell^2(\Lambda) \rightarrow \ell^2(\Lambda)} \leq 1$.

Using the maps K_m , R can be written succinctly as

$$R = \sum_{m=1}^{\infty} K_m \otimes R_m, \quad (1.23)$$

with unconditional convergence in $\mathcal{L}(\ell^2(\Lambda; V), \ell^2(\Lambda; V^*))$. By Lemma 1.1,

$$\|\mathbf{R}\|_{\ell^2(\Lambda; V) \rightarrow \ell^2(\Lambda; V^*)} \leq \gamma < 1. \quad (1.24)$$

In particular, $\|\mathbf{A}\| \leq (1 + \gamma)$ and $\|\mathbf{A}^{-1}\| \leq (1 - \gamma)^{-1}$.

We also define the operator $\mathbf{D} := T_V^* \mathcal{D} T_V$. This is just the Riesz isomorphism from $\ell^2(\Lambda; V)$ to $\ell^2(\Lambda; V^*)$. By [Git11c, Proposition 2.10],

$$(1 - \gamma)\mathbf{D} \leq \mathbf{A} \leq (1 + \gamma)\mathbf{D} \quad \text{and} \quad \frac{1}{1 + \gamma}\mathbf{D}^{-1} \leq \mathbf{A}^{-1} \leq \frac{1}{1 - \gamma}\mathbf{D}^{-1}. \quad (1.25)$$

In particular, using $\mathbf{A} = \mathbf{A}\mathbf{A}^{-1}\mathbf{A}$, we have

$$\frac{1}{1 + \gamma}\mathbf{A}\mathbf{D}^{-1}\mathbf{A} \leq \mathbf{A} \leq \frac{1}{1 - \gamma}\mathbf{A}\mathbf{D}^{-1}\mathbf{A}. \quad (1.26)$$

1.5 Galerkin Projection

Let \mathcal{W} be a closed subspace of $L_\pi^2(\Gamma; V)$. The Galerkin solution $\bar{u} \in \mathcal{W}$ is defined through the linear variational problem

$$\int_\Gamma \langle \mathbf{A}(y)\bar{u}(y), w(y) \rangle \, d\pi(y) = \int_\Gamma \langle f(y), w(y) \rangle \, d\pi(y) \quad \forall w \in \mathcal{W}. \quad (1.27)$$

Existence, uniqueness and quasi-optimality of \bar{u} follow since \mathcal{A} induces an inner product on $L_\pi^2(\Gamma; V)$ that is equivalent to the standard inner product, see [Git11c, Proposition 1.5].

For all $v \in \Lambda$, let W_v be a finite dimensional subspace of V , such that $W_v \neq \{0\}$ for only finitely many $v \in \Lambda$. It is particularly useful to consider spaces \mathcal{W} of the form

$$\mathcal{W} := \sum_{v \in \Lambda} W_v P_v. \quad (1.28)$$

The Galerkin operator on such a space has a similar structure to (1.20), with $A_{v\mu}$ replaced by its representation on suitable subspaces W_v of V , see [Git11c, Section 2].

2 Approximation of the Residual

2.1 Adaptive Application of the Stochastic Operator

We construct a sequence of approximations of \mathbf{R} by truncating the series (1.23). For all $M \in \mathbb{N}$, let

$$\mathbf{R}_{[M]} := \sum_{m=1}^M \mathbf{K}_m \otimes \mathbf{R}_m, \quad (2.1)$$

and $\mathbf{R}_{[0]} := 0$. For all $M \in \mathbb{N}$, let $\bar{e}_{\mathbf{R}, M}$ be given such that

$$\|\mathbf{R} - \mathbf{R}_{[M]}\|_{\ell^2(\Lambda; V) \rightarrow \ell^2(\Lambda; V^*)} \leq \bar{e}_{\mathbf{R}, M}. \quad (2.2)$$

For example, these bounds can be chosen as

$$\bar{e}_{\mathbf{R},M} := \sum_{m=M+1}^{\infty} \|\mathbf{R}_m\|_{V \rightarrow V^*} . \quad (2.3)$$

We assume that $(\bar{e}_{\mathbf{R},M})_{M=0}^{\infty}$ is nonincreasing and converges to 0, and also that the sequence of differences $(\bar{e}_{\mathbf{R},M} - \bar{e}_{\mathbf{R},M+1})_{M=0}^{\infty}$ is nonincreasing.

We consider a partitioning of a vector $\mathbf{w} \in \ell^2(\Lambda)$ into $\mathbf{w}_{[p]} := \mathbf{w}|_{\Lambda_p}$, $p = 1, \dots, P$, for disjoint index sets $\Lambda_p \subset \Lambda$. This can be approximate in that $\mathbf{w}_{[1]} + \dots + \mathbf{w}_{[P]}$ only approximates \mathbf{w} in $\ell^2(\Lambda)$. We think of $\mathbf{w}_{[1]}$ as containing the largest elements of \mathbf{w} , $\mathbf{w}_{[2]}$ the next largest, and so on.

Such a partitioning can be constructed by the approximate sorting algorithm

$$\text{BucketSort}[\mathbf{w}, \epsilon] \mapsto [(\mathbf{w}_{[p]})_{p=1}^P, (\Lambda_p)_{p=1}^P] , \quad (2.4)$$

which, given a finitely supported $\mathbf{w} \in \ell^2(\Lambda)$ and a threshold $\epsilon > 0$, returns index sets

$$\Lambda_p := \left\{ \mu \in \Lambda ; |v_{\mu}| \in (2^{-p/2} \|\mathbf{w}\|_{\ell^{\infty}}, 2^{-(p-1)/2} \|\mathbf{w}\|_{\ell^{\infty}}) \right\} \quad (2.5)$$

and $\mathbf{w}_{[p]} := \mathbf{w}|_{\Lambda_p}$, see [Met02, Bar05, GHS07, DSS09]. The integer P is minimal with

$$2^{-P/2} \|\mathbf{w}\|_{\ell^{\infty}(\Lambda)} \sqrt{\#\text{supp } \mathbf{w}} \leq \epsilon . \quad (2.6)$$

By [GHS07, Rem. 2.3] or [DSS09, Prop. 4.4], the number of operations and storage locations required by a call of $\text{BucketSort}[\mathbf{w}, \epsilon]$ is bounded by

$$\#\text{supp } \mathbf{w} + \max(1, \lceil \log(\|\mathbf{w}\|_{\ell^{\infty}(\Lambda)} \sqrt{\#\text{supp } \mathbf{w}} / \epsilon) \rceil) . \quad (2.7)$$

This analysis uses that every w_{μ} , $\mu \in \Lambda$, can be mapped to p with $\mu \in \Lambda_p$ in constant time by evaluating

$$p := \left\lceil 1 + 2 \log_2 \left(\frac{\|\mathbf{w}\|_{\ell^{\infty}(\Lambda)}}{|w_{\mu}|} \right) \right\rceil . \quad (2.8)$$

Alternatively, any standard comparison-based sorting algorithm can be used to construct the partitioning of \mathbf{w} , albeit with an additional logarithmic factor in the complexity.

The routine $\text{Apply}_{\mathbf{R}}[\mathbf{v}, \epsilon]$ adaptively approximates $\mathbf{R}\mathbf{v}$ in three distinct steps. First, the elements of \mathbf{v} are grouped according to their norm. Elements smaller than a certain tolerance are discarded. This truncation of the vector \mathbf{v} produces an error of at most $\delta \leq \epsilon/2$.

Next, a greedy algorithm is used to assign to each segment $\mathbf{v}_{[p]}$ of \mathbf{v} an approximation $\mathbf{R}_{[M_p]}$ of \mathbf{R} . Starting with $\mathbf{R}_{[M_p]} = \mathbf{0}$ for all $p = 1, \dots, \ell$, these approximations are refined iteratively until an estimate of the error is smaller than $\epsilon - \delta$.

Finally, the operations determined by the previous two steps are performed. Each multiplication $\mathbf{R}_m \mathbf{v}_{\mu}$ is performed just once, and copied to the appropriate entries of \mathbf{z} .

$\text{Apply}_{\mathbf{R}}[v, \epsilon] \mapsto z$

$[\cdot, (\Lambda_p)_{p=1}^P] \leftarrow \text{BucketSort} \left[(\|v_\mu\|_V)_{\mu \in \Lambda}, \frac{\epsilon}{2\bar{e}_{\mathbf{R},0}} \right]$

for $p = 1, \dots, P$ **do** $v_{[p]} \leftarrow (v_\mu)_{\mu \in \Lambda_p}$

Compute the minimal $\ell \in \{0, 1, \dots, P\}$ s.t. $\delta := \bar{e}_{\mathbf{R},0} \left\| v - \sum_{p=1}^{\ell} v_{[p]} \right\|_{\ell^2(\Lambda; V)} \leq \frac{\epsilon}{2}$

for $p = 1, \dots, P$ **do** $M_p \leftarrow 0$

while $\sum_{p=1}^{\ell} \bar{e}_{\mathbf{R}, M_p} \|v_{[p]}\|_{\ell^2(\Lambda; V)} > \epsilon - \delta$ **do**

$q \leftarrow \text{argmax}_{p=1, \dots, \ell} (\bar{e}_{\mathbf{R}, M_p} - \bar{e}_{\mathbf{R}, M_{p+1}}) \|v_{[p]}\|_{\ell^2(\Lambda; V)} / \#\Lambda_p$

$M_q \leftarrow M_q + 1$

$z = (z_\nu)_{\nu \in \Lambda} \leftarrow \mathbf{0}$

for $p = 1, \dots, \ell$ **do**

forall $\mu \in \Lambda_p$ **do**

for $m = 1, \dots, M_p$ **do**

$w \leftarrow R_m v_\mu$

$z_{\mu+\epsilon_m} \leftarrow z_{\mu+\epsilon_m} + \beta_{\mu_m+1}^m w$

if $\mu_m \geq 1$ **then** $z_{\mu-\epsilon_m} \leftarrow z_{\mu-\epsilon_m} + \beta_{\mu_m}^m w$

if $\alpha_{\mu_m}^m \neq 0$ **then** $z_\mu \leftarrow z_\mu + \alpha_{\mu_m}^m w$

Proposition 2.1. For any finitely supported $v \in \ell^2(\Lambda; V)$ and any $\epsilon > 0$, $\text{Apply}_{\mathbf{R}}[v, \epsilon]$ produces a finitely supported $z \in \ell^2(\Lambda; V^*)$ with

$$\#\text{supp } z \leq 3 \sum_{p=1}^{\ell} M_p \#\Lambda_p \quad (2.9)$$

and

$$\|Rv - z\|_{\ell^2(\Lambda; V^*)} \leq \delta + \eta_{\mathbf{M}} \leq \epsilon, \quad \eta_{\mathbf{M}} := \sum_{p=1}^{\ell} \bar{e}_{\mathbf{R}, M_p} \|v_{[p]}\|_{\ell^2(\Lambda; V)}, \quad (2.10)$$

where M_p refers to the final value of this variable in the call of $\text{Apply}_{\mathbf{R}}$. The total number of products $R_m v_\mu$ computed in $\text{Apply}_{\mathbf{R}}[v, \epsilon]$ is $\sigma_{\mathbf{M}} := \sum_{p=1}^{\ell} M_p \#\Lambda_p$. Furthermore, the vector $\mathbf{M} = (M_p)_{p=1}^{\ell}$ is optimal in the sense that if $\mathbf{N} = (N_p)_{p=1}^{\ell}$ with $\sigma_{\mathbf{N}} \leq \sigma_{\mathbf{M}}$ then $\eta_{\mathbf{N}} \geq \eta_{\mathbf{M}}$, and if $\eta_{\mathbf{N}} \leq \eta_{\mathbf{M}}$, then $\sigma_{\mathbf{N}} \geq \sigma_{\mathbf{M}}$.

Proof. The estimate (2.9) follows from the fact that each K_m has at most three nonzero entries per column, see (1.22). Since $\|\mathbf{R}\|_{\ell^2(\Lambda; V) \rightarrow \ell^2(\Lambda; V^*)} \leq \bar{e}_{\mathbf{R},0}$,

$$\left\| Rv - R \sum_{p=1}^{\ell} v_{[p]} \right\|_{\ell^2(\Lambda; V^*)} \leq \bar{e}_{\mathbf{R},0} \left\| v - \sum_{p=1}^{\ell} v_{[p]} \right\|_{\ell^2(\Lambda; V)} = \delta \leq \frac{\epsilon}{2}.$$

Due to (2.2) and the termination criterion in the greedy subroutine of `ApplyR`,

$$\sum_{p=1}^{\ell} \|\mathbf{R}v_{[p]} - \mathbf{R}_{[M_p]}v_{[p]}\|_{\ell^2(\Lambda; V^*)} \leq \sum_{p=1}^{\ell} \bar{\epsilon}_{\mathbf{R}, M_p} \|v_{[p]}\|_{\ell^2(\Lambda; V)} \leq \epsilon - \delta.$$

For the optimality property of the greedy algorithm, we refer to the more general statement [Git11a, Theorem 4.1.5]. \square

2.2 Computation of the Residual

We assume a solver for D is available such that for any $g \in V^*$ and any $\epsilon > 0$,

$$\text{Solve}_D[g, \epsilon] \mapsto v, \quad \|v - D^{-1}g\|_V \leq \epsilon. \quad (2.11)$$

For example, `SolveD` could be an adaptive wavelet method, see *e.g.* [CDD01, CDD02, GHS07], an adaptive frame method, see *e.g.* [Ste03, DFR07, DRW⁺07], or a finite element method with a posteriori error estimation, see *e.g.* [Dör96, MNS00, BDD04].

Furthermore, we assume that a routine

$$\text{RHS}_f[\epsilon] \mapsto \tilde{f} \quad (2.12)$$

is available to compute approximations $\tilde{f} = (\tilde{f}_v)_{v \in \Lambda}$ of f with $\#\text{supp } \tilde{f} < \infty$ and

$$\|f - \tilde{f}\|_{\ell^2(\Lambda; V^*)} \leq \epsilon \quad (2.13)$$

for any $\epsilon > 0$.

The routine `ResidualA, f` approximates the residual $f - Av$ up to a prescribed relative tolerance.

$$\text{Residual}_{A, f}[\epsilon, v, \eta_0, \chi, \omega, \alpha, \beta] \mapsto [w, \eta, \zeta]$$

$$\zeta \leftarrow \chi \eta_0$$

repeat

$$\left[\begin{array}{l} h = (h_v)_{v \in \Lambda} \leftarrow \text{RHS}_f[\beta(1 - \alpha)\zeta] - \text{Apply}_R[v, (1 - \beta)(1 - \alpha)\zeta] \\ w = (w_v)_{v \in \Lambda} \leftarrow (\text{Solve}_D[h_v, \alpha\zeta(\#\text{supp } h)^{-1/2}])_{v \in \Lambda} \\ \eta \leftarrow \|w - v\|_{\ell^2(\Lambda; V)} \\ \text{if } \zeta \leq \omega\eta \text{ or } \eta + \zeta \leq \epsilon \text{ then break} \\ \zeta \leftarrow \omega \frac{1 - \omega}{1 + \omega} (\eta + \zeta) \end{array} \right.$$

Proposition 2.2. *For any finitely supported $v = (v_v)_{v \in \Lambda} \in \ell^2(\Lambda; V)$, $\epsilon > 0$, $\eta_0 \geq 0$, $\chi > 0$, $\omega > 0$, $0 < \alpha < 1$ and $0 < \beta < 1$, a call of `ResidualA, f` computes $w \in \ell^2(\Lambda; V)$, $\eta \geq 0$ and $\zeta \geq 0$ with*

$$|\eta - \|r\|_{\ell^2(\Lambda; V^*)}| \leq \|w - v - D^{-1}r\|_{\ell^2(\Lambda; V)} = \|w - D^{-1}(f - Av)\|_{\ell^2(\Lambda; V)} \leq \zeta, \quad (2.14)$$

where $r = (r_v)_{v \in \Lambda} \in \ell^2(\Lambda; V^*)$ is the residual $r = f - Av$, and ζ satisfies either $\zeta \leq \omega\eta$ or $\eta + \zeta \leq \epsilon$.

Proof. By construction,

$$\|\mathbf{h} - (\mathbf{f} - \mathbf{R}\mathbf{v})\|_{\ell^2(\Lambda; V^*)} \leq \|\mathbf{h} - (\mathbf{f} - \mathbf{R}\mathbf{v})\|_{\ell^2(\Lambda; V^*)} \leq (1 - \alpha)\zeta.$$

Furthermore, using $\|\mathbf{w} - \mathbf{D}^{-1}\mathbf{h}\|_{\ell^2(\Lambda; V)} \leq \alpha\zeta$,

$$\|\mathbf{w} - \mathbf{D}^{-1}(\mathbf{f} - \mathbf{R}\mathbf{v})\|_{\ell^2(\Lambda; V)} \leq \|\mathbf{w} - \mathbf{D}^{-1}\mathbf{h}\|_{\ell^2(\Lambda; V)} + \|\mathbf{h} - (\mathbf{f} - \mathbf{R}\mathbf{v})\|_{\ell^2(\Lambda; V^*)} \leq \zeta.$$

The rest of (2.14) follows by triangle inequality with $\|\mathbf{r}\|_{\ell^2(\Lambda; V^*)} = \|\mathbf{D}^{-1}\mathbf{r}\|_{\ell^2(\Lambda; V)}$. \square

Remark 2.3. The tolerance ζ in $\text{Residual}_{A,f}$ is initialized as the product of an initial estimate η_0 of the residual and a parameter χ . The update

$$\zeta \leftarrow \omega \frac{1 - \omega}{1 + \omega} (\eta + \zeta) =: \zeta_1 \quad (2.15)$$

ensures a geometric decrease of ζ since if $\zeta > \omega\eta$, then

$$\zeta_1 = \omega \frac{1 - \omega}{1 + \omega} (\eta + \zeta) < \frac{1 - \omega}{1 + \omega} (\zeta + \omega\zeta) = (1 - \omega)\zeta. \quad (2.16)$$

Therefore, the total computational cost of the routine is proportional to that of the final iteration of the loop. Furthermore, if $\zeta > \omega\eta$, then also

$$\zeta_1 = \omega \frac{1 - \omega}{1 + \omega} (\eta + \zeta) > \omega(1 - \omega)\eta > \omega(\eta - \zeta). \quad (2.17)$$

The term $\eta - \zeta$ in the last expression of (2.17) is a lower bound for the true residual $\|\mathbf{r}\|_{\ell^2(\Lambda; V_D^*)}$. In this sense, the prescription (2.15) does not select an unnecessarily small tolerance.

Finally, if $\zeta \leq 2\omega(1 - \omega)^{-1}\eta$, then $\zeta_1 \leq \omega\eta$. If the next value of η is greater than or equal to the current value, this ensures that the termination criterion is met in the next iteration. For example, under the mild condition $\zeta \leq (1 + 4\omega - \omega^2)(1 - \omega)^{-2}\eta$, we have $\zeta_1 \leq 2\omega(1 - \omega)^{-1}\eta$. The loop can therefore be expected to terminate within three iterations. \lrcorner

Remark 2.4. In $\text{Residual}_{A,f}$, the tolerances of Solve_D are chosen such that the error tolerance $\alpha\zeta$ is equidistributed among all the nonzero indices of \mathbf{w} . This property is not required anywhere; Proposition 2.2 only uses that the total error in the computation of $\mathbf{D}^{-1}\mathbf{h}$ is no more than $\alpha\zeta$. Indeed, other strategies for selecting tolerances, e.g. based on additional a priori information, may be more efficient. Equidistributing the error among all the indices is a simple, practical starting point. \lrcorner

3 An Adaptive Solver

3.1 Refinement Strategy

We use the approximation of the residual described in Section 2 to refine a Galerkin subspace $\mathcal{W} \subset L^2_\pi(\Gamma; V)$ of the form (1.28). For some approximate solution v with $T_V v \in \mathcal{W}$, let w be the approximation of $D^{-1}(f - Rv)$ computed by $\text{Residual}_{A,f}$. We construct a space

$$\bar{\mathcal{W}} := \sum_{\mu \in \Lambda} \bar{W}_\mu P_\mu \supset \mathcal{W}, \quad (3.1)$$

with $\bar{W}_\mu \subset V$ finite dimensional, such that w can be approximated sufficiently in $\bar{\mathcal{W}}$. A simple choice is $\bar{W}_\mu := W_\mu + \text{span } w_\mu$, where $\mathcal{W} = \sum_\mu W_\mu P_\mu$.

We consider a multilevel setting. For each $\mu \in \text{supp } w \subset \Lambda$, let $W_\mu =: W_\mu^0 \subset W_\mu^1 \subset \dots$ be a scale of finite dimensional subspaces of V such that $\bigcup_{i=0}^\infty W_\mu^i$ is dense in V . To each space, we associate a cost $\dim W_\mu^i$ and an error $\|w_\mu - \Pi_\mu^i w_\mu\|_V^2$, where Π_μ^i denotes the orthogonal projection in V onto W_μ^i . In the construction of $\bar{\mathcal{W}}$, we use a greedy algorithm to minimize the dimension of $\bar{\mathcal{W}}$ under a constraint on the approximation error of w .

$\text{Refine}_D[\mathcal{W}, w, \epsilon] \mapsto [\bar{\mathcal{W}}, \bar{w}, \varrho]$

forall $\mu \in \text{supp } w$ **do** $j_\mu \leftarrow 0$

while $\sum_{\mu \in \text{supp } w} \|w_\mu - \Pi_\mu^{j_\mu} w_\mu\|_V^2 > \epsilon^2$ **do**

$v \leftarrow \underset{\mu \in \text{supp } w}{\text{argmax}} \frac{\|\Pi_\mu^{j_\mu+1} w_\mu - \Pi_\mu^{j_\mu} w_\mu\|_V^2}{\dim(W_\mu^{j_\mu+1} \setminus W_\mu^{j_\mu})}$
 $j_\mu \leftarrow j_\mu + 1$

forall $\mu \in \text{supp } w$ **do**

$\bar{W}_\mu \leftarrow W_\mu^{j_\mu}$
 $\bar{w}_\mu \leftarrow \Pi_\mu^{j_\mu} w_\mu$

$\varrho \leftarrow \left(\sum_{\mu \in \text{supp } w} \|w_\mu - \bar{w}_\mu\|_V^2 \right)^{1/2}$

Proposition 3.1. *If for every $\mu \in \text{supp } w$,*

$$\frac{\|\Pi_\mu^{i+1} w_\mu - \Pi_\mu^i w_\mu\|_V^2}{\dim(W_\mu^{i+1} \setminus W_\mu^i)} \geq \frac{\|\Pi_\mu^{j+1} w_\mu - \Pi_\mu^j w_\mu\|_V^2}{\dim(W_\mu^{j+1} \setminus W_\mu^j)} \quad \forall i \leq j, \quad (3.2)$$

then for any $\epsilon \geq 0$, a call of $\text{Refine}_D[\mathcal{W}, w, \epsilon]$ constructs a space $\bar{\mathcal{W}}$ of the form (3.1) and $T_V \bar{w} \in \bar{\mathcal{W}}$ satisfying

$$\varrho = \|w - \bar{w}\|_{\ell^2(\Lambda; V)} \leq \epsilon. \quad (3.3)$$

Furthermore, $\bar{\mathcal{W}}$ is minimal among all spaces of the form (3.1) with $\bar{W}_\mu = W_\mu^i$ and satisfying (3.3).

Proof. Equation (3.3) follows from the termination criterion in Refine_D . Convergence is ensured by (3.2) and $W_\mu^i \uparrow V$ for all μ . For the optimality property of the greedy algorithm, we refer to the more general statement [Git11a, Theorem 4.1.5]. \square

3.2 Adaptive Galerkin Method

Let $\|\cdot\|_A$ denote the energy norm on $\ell^2(\Lambda; V)$, i.e. $\|v\|_A := \sqrt{\langle Av, v \rangle}$. We assume that a routine

$$\text{Galerkin}_{A,f}[\mathcal{W}, \tilde{u}_0, \epsilon] \mapsto [\tilde{u}, \tau] \quad (3.4)$$

is available which, given a finite dimensional subspace \mathcal{W} of $L_\pi^2(\Gamma; V)$ of the form (1.28), and starting from the initial approximation \tilde{u}_0 , iteratively computes $\tilde{u} \in \ell^2(\Lambda; V)$ with $T_V \tilde{u} \in \mathcal{W}$ and

$$\|\tilde{u} - \bar{u}\|_A \leq \tau \leq \epsilon, \quad (3.5)$$

where $T_V \bar{u}$ is the Galerkin projection of u onto \mathcal{W} . An example of such a routine, based on a preconditioned conjugate gradient iteration, is given in [Git11c].

We combine the method $\text{Residual}_{A,f}$ for approximating the residual, Refine_D for refining the Galerkin subspace and $\text{Galerkin}_{A,f}$ for approximating the Galerkin projection, to an adaptive solver $\text{SolveGalerkin}_{A,f}$ similar to [CDD01, GHS07, DSS09].

$$\text{SolveGalerkin}_{A,f}[\epsilon, \gamma, \chi, \vartheta, \omega, \sigma, \alpha, \beta] \mapsto \mathbf{u}_\epsilon$$

$$\begin{aligned} & \mathcal{W}^{(0)} \leftarrow \{0\} \\ & \tilde{u}^{(0)} \leftarrow \mathbf{0} \\ & \delta_0 \leftarrow \sqrt{(1-\gamma)^{-1}} \|f\|_{\ell^2(\Lambda; V^*)} \\ & \mathbf{for } k = 0, 1, 2, \dots \mathbf{ do} \\ & \quad \left[\begin{aligned} & [\mathbf{w}_k, \eta_k, \zeta_k] \leftarrow \text{Residual}_{A,f}[\epsilon \sqrt{1-\gamma}, \tilde{u}^{(k)}, \delta_k, \chi, \omega, \alpha, \beta] \\ & \bar{\delta}_k \leftarrow (\eta_k + \zeta_k) / \sqrt{1-\gamma} \\ & \mathbf{if } \min(\delta_k, \bar{\delta}_k) \leq \epsilon \mathbf{ then break} \\ & [\mathcal{W}^{(k+1)}, \bar{\mathbf{w}}_k, \varrho_k] \leftarrow \text{Refine}_D[\mathcal{W}^{(k)}, \mathbf{w}_k, \sqrt{\eta_k^2 - (\zeta_k + \vartheta(\eta_k + \zeta_k))^2}] \\ & \bar{\vartheta}_k \leftarrow (\sqrt{\eta_k^2 - \varrho_k^2} - \zeta_k) / (\eta_k + \zeta_k) \\ & [\tilde{u}^{(k+1)}, \tau_{k+1}] \leftarrow \text{Galerkin}_{A,f}[\mathcal{W}^{(k+1)}, \bar{\mathbf{w}}_k, \sigma \min(\delta_k, \bar{\delta}_k)] \\ & \delta_{k+1} \leftarrow \tau_{k+1} + \sqrt{1 - \bar{\vartheta}_k^2 (1-\gamma)(1+\gamma)^{-1}} \min(\delta_k, \bar{\delta}_k) \end{aligned} \right. \\ & \mathbf{u}_\epsilon \leftarrow \tilde{u}^{(k)} \end{aligned}$$

3.3 Convergence of the Adaptive Solver

The convergence analysis of $\text{SolveGalerkin}_{A,f}$ is based [CDD01, Lemma 4.1], which generalizes to our vector setting for Galerkin spaces \mathcal{W} of the form (1.28). Let $\Pi_{\mathcal{W}}$

denote the orthogonal projection in $\ell^2(\Lambda; V)$ onto $T_V^{-1}\mathcal{W}$, and let $\hat{\Pi}_{\mathcal{W}} := D\Pi_{\mathcal{W}}D^{-1}$ be the orthogonal projection in $\ell^2(\Lambda; V^*)$ onto $DT_V^{-1}\mathcal{W} = T_V^*\mathcal{D}\mathcal{W}$.

Proposition 3.2. *Let \mathcal{W} be as in (1.28), and $\vartheta \in [0, 1]$. Let $v \in \mathcal{W}$ with*

$$\|\hat{\Pi}_{\mathcal{W}}(f - Av)\|_{\ell^2(\Lambda; V^*)} \geq \vartheta \|f - Av\|_{\ell^2(\Lambda; V^*)}. \quad (3.6)$$

Then the Galerkin projection \bar{u} of u onto \mathcal{W} satisfies

$$\|u - \bar{u}\|_A \leq \sqrt{1 - \vartheta^2 \frac{1 - \gamma}{1 + \gamma}} \|u - v\|_A. \quad (3.7)$$

Proof. Due to (3.6),

$$\begin{aligned} \|\bar{u} - v\|_A &\geq \|A\|^{-1/2} \|A(\bar{u} - v)\|_{\ell^2(\Lambda; V^*)} \geq \|A\|^{-1/2} \|\hat{\Pi}_{\mathcal{W}}(f - Av)\|_{\ell^2(\Lambda; V^*)} \\ &\geq \|A\|^{-1/2} \vartheta \|f - Av\|_{\ell^2(\Lambda; V^*)} \geq \|A\|^{-1/2} \|A^{-1}\|^{-1/2} \vartheta \|u - v\|_A. \end{aligned}$$

By Galerkin orthogonality,

$$\|u - \bar{u}\|_A^2 = \|u - v\|_A^2 - \|\bar{u} - v\|_A^2 \leq (1 - \vartheta^2 \|A\|^{-1} \|A^{-1}\|^{-1}) \|u - v\|_A^2.$$

The assertion follows using the estimates $\|A\| \leq (1 + \gamma)$ and $\|A^{-1}\| \leq (1 - \gamma)^{-1}$, which follow from (1.24). \square

Lemma 3.3. *If $\vartheta > 0$, $\omega > 0$, and $\omega + \vartheta + \omega\vartheta \leq 1$, then the space $\mathcal{W}^{(k+1)}$ in $\text{SolveGalerkin}_{A,f}$ is such that*

$$\|\hat{\Pi}_{\mathcal{W}^{(k+1)}} r_k\|_{\ell^2(\Lambda; V^*)} \geq \bar{\vartheta}_k \|r_k\|_{\ell^2(\Lambda; V^*)} \quad (3.8)$$

where $r_k := f - A\tilde{u}^{(k)}$ is the residual at iteration $k \in \mathbb{N}_0$, and $\bar{\vartheta}_k \geq \vartheta$.

Proof. We abbreviate $z := w_k - \tilde{u}^{(k)}$. Due to $\zeta_k \leq \omega\eta_k$, the assumption $\omega + \vartheta + \omega\vartheta \leq 1$ implies $\zeta_k + \vartheta(\eta_k + \zeta_k) \leq \eta_k$. Thus the tolerance in Refine_D is nonnegative. Since $\tilde{u}^{(k)} \in \mathcal{W}^{(k)} \subset \mathcal{W}^{(k+1)}$, Proposition 3.1 implies

$$\varrho_k = \|w_k - \bar{w}_k\|_{\ell^2(\Lambda; V)} = \|w_k - \Pi_{\mathcal{W}^{(k+1)}} w_k\|_{\ell^2(\Lambda; V)} = \|z - \Pi_{\mathcal{W}^{(k+1)}} z\|_{\ell^2(\Lambda; V)}.$$

Consequently,

$$\|\Pi_{\mathcal{W}^{(k+1)}} z\|_{\ell^2(\Lambda; V)}^2 = \|z\|_{\ell^2(\Lambda; V)}^2 - \|z - \Pi_{\mathcal{W}^{(k+1)}} z\|_{\ell^2(\Lambda; V)}^2 = \eta_k^2 - \varrho_k^2.$$

Furthermore, since $\Pi_{\mathcal{W}^{(k+1)}}$ has norm one, Proposition 2.2 implies

$$\begin{aligned} \|\Pi_{\mathcal{W}^{(k+1)}} z\|_{\ell^2(\Lambda; V)} - \|\hat{\Pi}_{\mathcal{W}^{(k+1)}} r_k\|_{\ell^2(\Lambda; V^*)} &\leq \|\Pi_{\mathcal{W}^{(k+1)}}(z - D^{-1}r_k)\|_{\ell^2(\Lambda; V)} \\ &\leq \|z - D^{-1}r_k\|_{\ell^2(\Lambda; V)} \leq \zeta_k. \end{aligned}$$

Combining these estimates, we have

$$\|\hat{\Pi}_{\mathcal{W}^{(k+1)}} \mathbf{r}_k\|_{\ell^2(\Lambda; V^*)} \geq \|\Pi_{\mathcal{W}^{(k+1)}} \mathbf{z}\|_{\ell^2(\Lambda; V)} - \zeta_k = \sqrt{\eta_k^2 - \varrho_k^2} - \zeta_k,$$

and (3.8) follows using $\|\mathbf{r}_k\|_{\ell^2(\Lambda; V^*)} \leq \eta_k + \zeta_k$. Finally, $\varrho_k^2 \leq \eta_k^2 - (\zeta_k + \vartheta(\eta_k + \zeta_k))^2$ implies $\sqrt{\eta_k^2 - \varrho_k^2} \geq \zeta_k + \vartheta(\eta_k + \zeta_k)$, and therefore $\bar{\delta}_k = (\sqrt{\eta_k^2 - \varrho_k^2} - \zeta_k)/(\eta_k + \zeta_k) \geq \vartheta$. \square

Theorem 3.4. *If $\epsilon > 0$, $\chi > 0$, $\vartheta > 0$, $\omega > 0$, $\omega + \vartheta + \omega\vartheta \leq 1$, $0 < \alpha < 1$, $0 < \beta < 1$ and $0 < \sigma < 1 - \sqrt{1 - \vartheta^2(1 - \gamma)(1 + \gamma)^{-1}}$, then $\text{SolveGalerkin}_{\mathbf{A}, \mathbf{f}}[\epsilon, \gamma, \chi, \vartheta, \omega, \sigma, \alpha, \beta]$ constructs a finitely supported $\mathbf{u}_\epsilon \in \ell^2(\Lambda; V)$ with*

$$\|\mathbf{u} - \mathbf{u}_\epsilon\|_{\mathbf{A}} \leq \epsilon. \quad (3.9)$$

Moreover,

$$\sqrt{\frac{1 - \gamma}{1 + \gamma} \frac{1 - \omega}{1 + \omega}} \bar{\delta}_k \leq \|\mathbf{u} - \tilde{\mathbf{u}}^{(k)}\|_{\mathbf{A}} \leq \min(\delta_k, \bar{\delta}_k) \quad (3.10)$$

for all $k \in \mathbb{N}_0$ reached by $\text{SolveGalerkin}_{\mathbf{A}, \mathbf{f}}$.

Proof. Due to the termination criterion of $\text{SolveGalerkin}_{\mathbf{A}, \mathbf{f}}$, it suffices to show (3.10).

For $k = 0$, since $\|\mathbf{u}\|_{\ell^2(\Lambda; V)} \leq \|\mathbf{A}^{-1}\|^{1/2} \|\mathbf{u}\|_{\mathbf{A}}$,

$$\|\mathbf{u} - \tilde{\mathbf{u}}^{(0)}\|_{\mathbf{A}}^2 = \|\mathbf{u}\|_{\mathbf{A}}^2 = \langle \mathbf{f}, \mathbf{u} \rangle_{\ell^2(\Lambda; V)} \leq \|\mathbf{f}\|_{\ell^2(\Lambda; V^*)} \|\mathbf{u}\|_{\ell^2(\Lambda; V)} \leq \delta_0 \|\mathbf{u}\|_{\mathbf{A}}.$$

Let $\|\mathbf{u} - \tilde{\mathbf{u}}^{(k)}\|_{\mathbf{A}} \leq \delta_k$ for some $k \in \mathbb{N}_0$. Abbreviating $\mathbf{r}_k := \mathbf{f} - \mathbf{A}\tilde{\mathbf{u}}^{(k)}$, using (1.26) then (2.14), we have

$$\|\mathbf{u} - \tilde{\mathbf{u}}^{(k)}\|_{\mathbf{A}} \leq \frac{1}{\sqrt{1 - \gamma}} \|\mathbf{r}_k\|_{\ell^2(\Lambda; V^*)} \leq \frac{\zeta_k + \eta_k}{\sqrt{1 - \gamma}} = \bar{\delta}_k.$$

If $\min(\delta_k, \bar{\delta}_k) > \epsilon$, then $\zeta_k \leq \omega\eta_k$ by Proposition 2.2. Due to Lemma 3.3, Proposition 3.2 implies

$$\|\mathbf{u} - \bar{\mathbf{u}}\|_{\mathbf{A}} \leq \sqrt{1 - \vartheta_k^2 \frac{1 - \gamma}{1 + \gamma}} \min(\delta_k, \bar{\delta}_k),$$

where $\bar{\mathbf{u}}$ is the exact Galerkin projection of \mathbf{u} onto $\mathcal{W}^{(k+1)}$. By (3.5), $\tilde{\mathbf{u}}^{(k+1)}$ approximates $\bar{\mathbf{u}}$ up to an error of at most $\tau_{k+1} \leq \sigma \min(\delta_k, \bar{\delta}_k)$ in the norm $\|\cdot\|_{\mathbf{A}}$. It follows by triangle inequality that $\|\mathbf{u} - \tilde{\mathbf{u}}^{(k+1)}\|_{\mathbf{A}} \leq \delta_{k+1}$.

To show the other inequality in (3.10), we note that for any $k \in \mathbb{N}_0$,

$$\|\mathbf{u} - \tilde{\mathbf{u}}^{(k)}\|_{\mathbf{A}} \geq \frac{1}{\sqrt{1 + \gamma}} \|\mathbf{r}_k\|_{\ell^2(\Lambda; V^*)} \geq \frac{\eta_k - \zeta_k}{\sqrt{1 + \gamma}} = \sqrt{\frac{1 - \gamma}{1 + \gamma} \frac{\eta_k - \zeta_k}{\eta_k + \zeta_k}} \bar{\delta}_k,$$

and $(\eta_k - \zeta_k)(\eta_k + \zeta_k)^{-1} \geq (1 - \omega)(1 + \omega)^{-1}$.

Finally, since

$$\delta_k \leq \left(\sigma + \sqrt{1 - \vartheta^2(1 - \gamma)(1 + \gamma)^{-1}} \right)^k \delta_0$$

and $\sigma + \sqrt{1 - \vartheta^2(1 - \gamma)(1 + \gamma)^{-1}} < 1$ by assumption, the iteration does terminate. \square

4 Optimality Properties

4.1 A Semidiscrete Algorithm

We consider the adaptive method from Section 3 with no discretization in V , *i.e.* with Galerkin subspaces of the form $\ell^2(\Xi; V) \subset \ell^2(\Lambda; V)$ for a finite subset $\Xi \subset \Lambda$. Formally replacing $\mathcal{W}^{(k)}$ by a set of active indices $\Xi^{(k)}$, $\text{SolveGalerkin}_{\mathbf{A},f}$ naturally extends to this setting.

In the subroutine $\text{Residual}_{\mathbf{A},f}$, we assume that Solve_D inverts D exactly. The parameter α can thus be set to zero.

In the subsequent refinement step, $\Xi^{(k)}$ is augmented by sufficiently many elements of $\text{supp } w_k$ to represent w_k to the desired accuracy. The method Refine_D reduces to ordering $\text{supp } w_k$ according to $\|w_{k,v}\|_V$ and selecting the most important contributions.

In $\text{Galerkin}_{\mathbf{A},f}$, we assume that operations in V can be performed exactly, and that the Galerkin projection of \mathbf{u} onto $\ell^2(\Xi^{(k+1)}; V)$ can be approximated *e.g.* by a conjugate gradient iteration.

4.2 Optimal Choice of Subspaces

For $v \in \ell^2(\Lambda; V)$ and $N \in \mathbb{N}_0$, let $P_N(v)$ be a best N -term approximation of v , that is, $P_N(v)$ is an element of $\ell^2(\Lambda; V)$ that minimizes $\|v - v_N\|_{\ell^2(\Lambda; V)}$ over $v_N \in \ell^2(\Lambda; V)$ with $\#\text{supp } v_N \leq N$. For $s \in (0, \infty)$, we define

$$\|v\|_{\mathcal{A}^s(\Lambda; V)} := \sup_{N \in \mathbb{N}_0} (N+1)^s \|v - P_N(v)\|_{\ell^2(\Lambda; V)} \quad (4.1)$$

and

$$\mathcal{A}^s(\Lambda; V) := \left\{ v \in \ell^2(\Lambda; V) ; \|v\|_{\mathcal{A}^s(\Lambda; V)} < \infty \right\}. \quad (4.2)$$

By definition, an optimal approximation in $\ell^2(\Lambda; V)$ of $v \in \mathcal{A}^s(\Lambda; V)$ with error tolerance $\epsilon > 0$ consists of $\mathcal{O}(\epsilon^{-1/s})$ nonzero coefficients in V .

For any $\Xi \subset \Lambda$, let Π_Ξ denote the orthogonal projection in $\ell^2(\Lambda; V^*)$ onto $\ell^2(\Xi; V^*)$. The following statement is adapted from [GHS07, Lemma 2.1] and [DSS09, Lemma 4.1].

Lemma 4.1. *Let $\Xi^{(0)}$ be a finite subset of Λ and $v \in \ell^2(\Xi^{(0)}; V)$. If*

$$0 \leq \hat{\delta} < \sqrt{\frac{1-\gamma}{1+\gamma}} \quad (4.3)$$

and $\Xi^{(0)} \subset \Xi^{(1)} \subset \Lambda$ with

$$\#\Xi^{(1)} \leq \bar{c} \min \left\{ \#\Xi ; \Xi^{(0)} \subset \Xi, \left\| \Pi_\Xi (f - \mathbf{A}v) \right\|_{\ell^2(\Lambda; V^*)} \geq \hat{\delta} \|f - \mathbf{A}v\|_{\ell^2(\Lambda; V^*)} \right\} \quad (4.4)$$

for a $\bar{c} \geq 1$, then

$$\#(\Xi^{(1)} \setminus \Xi^{(0)}) \leq \bar{c} \min \left\{ \#\hat{\Xi} ; \hat{\Xi} \subset \Lambda, \|\mathbf{u} - \hat{\mathbf{u}}\|_{\mathbf{A}} \leq \tau \|\mathbf{u} - v\|_{\mathbf{A}} \right\} \quad (4.5)$$

for $\tau = \sqrt{1 - \hat{\delta}^2(1+\gamma)(1-\gamma)^{-1}}$, where $\hat{\mathbf{u}}$ denotes the Galerkin projection of \mathbf{u} onto $\ell^2(\hat{\Xi}; V)$.

Proof. Let $\hat{\Xi}$ be as in (4.5) and $\check{\Xi} := \Xi^{(0)} \cup \hat{\Xi}$. Furthermore, let $\hat{\mathbf{u}}$ and $\check{\mathbf{u}}$ denote the Galerkin solutions in $\ell^2(\hat{\Xi}; V)$ and $\ell^2(\check{\Xi}; V)$, respectively. Since $\hat{\Xi} \subset \check{\Xi}$, $\|\mathbf{u} - \check{\mathbf{u}}\|_{\mathbf{A}} \leq \|\mathbf{u} - \hat{\mathbf{u}}\|_{\mathbf{A}}$, and by Galerkin orthogonality,

$$\|\check{\mathbf{u}} - \mathbf{v}\|_{\mathbf{A}}^2 = \|\mathbf{u} - \mathbf{v}\|_{\mathbf{A}}^2 - \|\mathbf{u} - \check{\mathbf{u}}\|_{\mathbf{A}}^2 \geq (1 - \tau^2) \|\mathbf{u} - \mathbf{v}\|_{\mathbf{A}}^2 = \hat{\vartheta}^2(1 + \gamma)(1 - \gamma^{-1}) \|\mathbf{u} - \mathbf{v}\|_{\mathbf{A}}^2 .$$

Therefore, using $\kappa(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \leq (1 + \gamma)(1 - \gamma)^{-1}$,

$$\begin{aligned} \|\Pi_{\check{\Xi}}(\mathbf{f} - \mathbf{A}\mathbf{v})\|_{\ell^2(\Lambda; V^*)} &= \|\mathbf{A}(\check{\mathbf{u}} - \mathbf{v})\|_{\ell^2(\Lambda; V^*)} \geq \|\mathbf{A}^{-1}\|^{-1/2} \|\check{\mathbf{u}} - \mathbf{v}\|_{\mathbf{A}} \\ &\geq \hat{\vartheta} \|\mathbf{A}\|^{1/2} \|\mathbf{u} - \mathbf{v}\|_{\mathbf{A}} \geq \hat{\vartheta} \|\mathbf{f} - \mathbf{A}\mathbf{v}\|_{\ell^2(\Lambda; V^*)} . \end{aligned}$$

By (4.4), $\#\Xi^{(1)} \leq \bar{c}\#\check{\Xi}$, and consequently

$$\#(\Xi^{(1)} \setminus \Xi^{(0)}) \leq \bar{c}\#(\check{\Xi} \setminus \Xi^{(0)}) \leq \bar{c}\#\hat{\Xi} . \quad \square$$

We use Lemma 4.1 to show that, under additional assumptions on the parameters, the index sets $\Xi^{(k)}$ generated by the semidiscrete version of `SolveGalerkinA,f` are of optimal size, up to a constant factor.

Theorem 4.2. *If the conditions of Theorem 3.4 are satisfied,*

$$\hat{\vartheta} := \frac{\vartheta(1 + \omega) + 2\omega}{1 - \omega} < \sqrt{\frac{1 - \gamma}{1 + \gamma}} , \quad (4.6)$$

and $\mathbf{u} \in \mathcal{F}^s(\Lambda; V)$ for an $s > 0$, then for all $k \in \mathbb{N}_0$ reached by `SolveGalerkinA,f`,

$$\#\Xi^{(k)} \leq 2 \frac{(\varrho/\tau)^{1/s}}{1 - \varrho^{1/s}} \left(\frac{(1 + \gamma)(1 + \omega)}{(1 - \gamma)(1 - \omega)} \right)^{1/s} \|\mathbf{u} - \tilde{\mathbf{u}}^{(k)}\|_{\ell^2(\Lambda; V)}^{-1/s} \|\mathbf{u}\|_{\mathcal{F}^s(\Lambda; V)}^{1/s} \quad (4.7)$$

with $\varrho = \sigma + \sqrt{1 - \vartheta^2(1 - \gamma)(1 + \gamma)^{-1}}$ and $\tau = \sqrt{1 - \hat{\vartheta}^2(1 + \gamma)(1 - \gamma)^{-1}}$.

Proof. Let $k \in \mathbb{N}_0$, $\mathbf{r}_k = \mathbf{f} - \mathbf{A}\tilde{\mathbf{u}}^{(k)}$. Also, let $\varrho = (\varrho_v)_{v \in \Lambda}$, $\varrho_v := \|w_{k,v}\|_V$ for the approximation $w_k = (w_{k,v})_{v \in \Lambda}$ of $\mathbf{D}^{-1}\mathbf{r}_k$ computed in `ResidualA,f`, and let $\Delta \subset \text{supp } w_k$ denote the active indices selected by `RefineD`.

We note that for $\alpha := \omega + \vartheta + \omega\vartheta$, we have $\vartheta = \frac{\alpha - \omega}{1 + \omega}$ and $\hat{\vartheta} = \frac{\alpha + \omega}{1 - \omega}$. Let $\Xi^{(k)} \subset \check{\Xi} \subset \Lambda$ satisfy $\|\Pi_{\check{\Xi}}\mathbf{r}_k\|_{\ell^2(\Lambda; V^*)} \geq \hat{\vartheta} \|\mathbf{r}_k\|_{\ell^2(\Lambda; V^*)}$. Then

$$\begin{aligned} \hat{\vartheta} \|\varrho\|_{\ell^2(\Lambda)} &\leq \hat{\vartheta} \|\mathbf{r}_k\|_{\ell^2(\Lambda; V^*)} + \hat{\vartheta}\omega \|\varrho\|_{\ell^2(\Lambda)} \\ &\leq \|\Pi_{\check{\Xi}}\mathbf{r}_k\|_{\ell^2(\Lambda; V^*)} + \hat{\vartheta}\omega \|\varrho\|_{\ell^2(\Lambda)} \leq \|\Pi_{\check{\Xi}}\varrho\|_{\ell^2(\Lambda)} + (1 + \hat{\vartheta})\omega \|\varrho\|_{\ell^2(\Lambda)} \end{aligned}$$

and since $\hat{\vartheta} - (1 + \hat{\vartheta})\omega = \alpha$, it follows that $\|\Pi_{\check{\Xi}}\varrho\|_{\ell^2(\Lambda)} \geq \alpha \|\varrho\|_{\ell^2(\Lambda)}$. By construction, Δ is a set of minimal cardinality with $\|\Pi_{\Delta}\varrho\|_{\ell^2(\Lambda)} \geq \bar{\alpha} \|\varrho\|_{\ell^2(\Lambda)}$ for $\bar{\alpha} := \zeta_k \eta_k^{-1} + \vartheta(1 + \zeta_k \eta_k^{-1}) \leq \alpha$.

Consequently, $\#(\Xi^{(k+1)} \setminus \Xi^{(k)}) \leq \#\Delta \leq \#\bar{\Xi}$. Since this holds for any $\bar{\Xi}$, using $\#\Xi^{(k)} \leq \bar{\Xi}$, it follows that

$$\#\Xi^{(k+1)} \leq 2 \min \left\{ \#\bar{\Xi} ; \Xi^{(k)} \subset \bar{\Xi} \subset \Lambda, \|\Pi_{\bar{\Xi}} \mathbf{r}_k\|_{\ell^2(\Lambda; V^*)} \geq \hat{\mathfrak{D}} \|\mathbf{r}_k\|_{\ell^2(\Lambda; V^*)} \right\}.$$

Lemma 4.1 implies

$$\#(\bar{\Xi}^{(k+1)} \setminus \bar{\Xi}^{(k)}) \leq 2 \min \left\{ \#\hat{\Xi} ; \hat{\Xi} \subset \Lambda, \|\mathbf{u} - \hat{\mathbf{u}}\|_{\mathbf{A}} \leq \tau \|\mathbf{u} - \tilde{\mathbf{u}}^{(k)}\|_{\mathbf{A}} \right\}$$

with $\tau = \sqrt{1 - \hat{\mathfrak{D}}^2(1 + \gamma)(1 - \gamma)^{-1}}$, where $\hat{\mathbf{u}}$ denotes the Galerkin projection of \mathbf{u} onto $\ell^2(\hat{\Xi}; V)$.

Let $N \in \mathbb{N}_0$ be maximal with $\|\mathbf{u} - P_N(\mathbf{u})\|_{\ell^2(\Lambda; V)} > \tau(1 + \gamma)^{-1/2} \|\mathbf{u} - \tilde{\mathbf{u}}^{(k)}\|_{\mathbf{A}}$, where $P_N(\mathbf{u})$ is a best N -term approximation of \mathbf{u} . By (4.1),

$$N + 1 \leq \|\mathbf{u} - P_N(\mathbf{u})\|_{\ell^2(\Lambda; V)}^{-1/s} \|\mathbf{u}\|_{\mathcal{A}^s(\Lambda; V)}^{1/s} \leq \tau^{-1/s} (1 + \gamma)^{1/2s} \|\mathbf{u} - \tilde{\mathbf{u}}^{(k)}\|_{\mathbf{A}}^{-1/s} \|\mathbf{u}\|_{\mathcal{A}^s(\Lambda; V)}^{1/s}.$$

For $\Xi_{N+1} := \text{supp } P_{N+1}(\mathbf{u})$, by maximality of N ,

$$\|\mathbf{u} - \bar{\mathbf{u}}_{N+1}\|_{\mathbf{A}} \leq \|\mathbf{u} - P_{N+1}(\mathbf{u})\|_{\mathbf{A}} \leq (1 + \gamma)^{1/2} \|\mathbf{u} - P_{N+1}(\mathbf{u})\|_{\ell^2(\Lambda; V)} \leq \tau \|\mathbf{u} - \tilde{\mathbf{u}}^{(k)}\|_{\mathbf{A}}$$

for the Galerkin solution $\bar{\mathbf{u}}_{N+1}$ in $\ell^2(\Xi_{N+1}; V)$, and thus

$$\#(\Xi^{(k+1)} \setminus \Xi^{(k)}) \leq 2(N + 1) \leq 2\tau^{-1/s} (1 + \gamma)^{1/2s} \|\mathbf{u} - \tilde{\mathbf{u}}^{(k)}\|_{\mathbf{A}}^{-1/s} \|\mathbf{u}\|_{\mathcal{A}^s(\Lambda; V)}^{1/s}.$$

Furthermore, by Theorem 3.4,

$$\|\mathbf{u} - \tilde{\mathbf{u}}^{(k)}\|_{\mathbf{A}}^{-1/s} \leq \left(\sqrt{\frac{1 - \gamma}{1 + \gamma}} \frac{1 - \omega}{1 + \omega} \bar{\delta}_k \right)^{-1/s}.$$

We estimate the cardinality of $\Xi^{(k)}$ by slicing it into increments and applying the above estimates,

$$\begin{aligned} \#\Xi^{(k)} &= \sum_{j=0}^{k-1} \#(\Xi^{(j+1)} \setminus \Xi^{(j)}) \leq 2\tau^{-1/s} (1 + \gamma)^{1/2s} \|\mathbf{u}\|_{\mathcal{A}^s(\Lambda; V)}^{1/s} \sum_{j=0}^{k-1} \|\mathbf{u} - \tilde{\mathbf{u}}^{(j)}\|_{\mathbf{A}}^{-1/s} \\ &\leq 2 \left(\frac{\tau(1 - \gamma)^{1/2}(1 - \omega)}{(1 + \gamma)(1 + \omega)} \right)^{-1/s} \|\mathbf{u}\|_{\mathcal{A}^s(\Lambda; V)}^{1/s} \sum_{j=0}^{k-1} \bar{\delta}_j^{-1/s}. \end{aligned}$$

By definition, $\delta_k \leq \varrho^{k-j} \bar{\delta}_j$. Therefore,

$$\sum_{j=0}^{k-1} \bar{\delta}_j^{-1/s} \leq \delta_k^{-1/s} \sum_{j=0}^{k-1} \varrho^{(k-j)/s} = \delta_k^{-1/s} \sum_{i=1}^k \varrho^{i/s} = \frac{\varrho^{1/s} \delta_k^{-1/s}}{1 - \varrho^{1/s}}.$$

The assertion follows using

$$(1 - \gamma)^{1/2} \|\mathbf{u} - \tilde{\mathbf{u}}^{(k)}\|_{\ell^2(\Lambda; V)} \leq \|\mathbf{u} - \tilde{\mathbf{u}}^{(k)}\|_{\mathbf{A}} \leq \delta_k. \quad \square$$

4.3 Complexity Estimate

We first cite a result due to Stechkin connecting the order of summability of a sequence to the convergence of best N -term approximations in a weaker sequence norm, see *e.g.* [CDS10b, DeV98]. Note that, although it is formulated only for nonnegative sequences, Lemma 4.3 applies directly to *e.g.* Lebesgue–Bochner spaces of Banach space valued sequences by passing to the norms of the elements of such sequences. Also, it applies to sequences with arbitrary countable index sets by choosing a decreasing rearrangement.

Lemma 4.3. *Let $0 < p \leq q$ and let $c = (c_n)_{n=1}^\infty \in \ell^2$ with $0 \leq c_{n+1} \leq c_n$ for all $n \in \mathbb{N}$. Then*

$$\left(\sum_{n=N+1}^{\infty} c_n^q \right)^{1/q} \leq (N+1)^{-r} \|c\|_{\ell^p}, \quad r := \frac{1}{p} - \frac{1}{q} \geq 0 \quad (4.8)$$

for all $N \in \mathbb{N}_0$.

Proof. Due to the elementary estimate

$$\|c\|_{\ell^p}^p = \sum_{i=1}^{\infty} c_i^p \geq \sum_{i=1}^n c_i^p \geq \sum_{i=1}^n c_n^p = n c_n^p,$$

we have $c_n \leq n^{-1/p} \|c\|_{\ell^p}$ for all $n \in \mathbb{N}$. Therefore, using $q - p \geq 0$,

$$\sum_{n=N+1}^{\infty} c_n^q \leq \sum_{n=N+1}^{\infty} c_n^p c_{N+1}^{q-p} \leq \|c\|_{\ell^p}^p (N+1)^{-(q-p)/p} \|c\|_{\ell^p}^{q-p} = (N+1)^{-r q} \|c\|_{\ell^p}^q$$

for all $N \in \mathbb{N}_0$, with r as in (4.8). □

Proposition 4.4. *Let $s > 0$. If either*

$$\|R_m\|_{V \rightarrow V^*} \leq s \delta_{\mathbf{R},s} (m+1)^{-s-1} \quad \forall m \in \mathbb{N} \quad (4.9)$$

or

$$\left(\sum_{m=1}^{\infty} \|R_m\|_{V \rightarrow V^*}^{\frac{1}{s+1}} \right)^{s+1} \leq \delta_{\mathbf{R},s}, \quad (4.10)$$

then

$$\|\mathbf{R} - \mathbf{R}_{[M]}\|_{\ell^2(\Lambda; V) \rightarrow \ell^2(\Lambda; V^*)} \leq \delta_{\mathbf{R},s} (M+1)^{-s} \quad \forall M \in \mathbb{N}_0. \quad (4.11)$$

Proof. By (1.23) and (2.1), using $\|K_m\|_{\ell^2(\Lambda) \rightarrow \ell^2(\Lambda)} \leq 1$,

$$\|\mathbf{R} - \mathbf{R}_{[M]}\|_{\ell^2(\Lambda; V) \rightarrow \ell^2(\Lambda; V^*)} \leq \sum_{m=M+1}^{\infty} \|R_m\|_{V \rightarrow V^*}.$$

If (4.9) holds, then (4.11) follows using

$$\sum_{m=M+1}^{\infty} (m+1)^{-s-1} \leq \int_{M+1}^{\infty} t^{-s-1} dt = \frac{1}{s} (M+1)^{-s}.$$

If (4.10) is satisfied, then

$$\sum_{m=M+1}^{\infty} \|R_m\|_{V \rightarrow V^*} \leq \left(\sum_{m=1}^{\infty} \|R_m\|_{V \rightarrow V^*}^{\frac{1}{s+1}} \right)^{s+1} (M+1)^{-s}$$

by Lemma 4.3. \square

Remark 4.5. If the assumptions of Proposition 4.4 are satisfied for all $s \in (0, s^*)$, then the operator R is s^* -compressible with sparse approximations $R_{[M]}$. In this case, R is a bounded linear map from $\mathcal{A}^s(\Lambda; V)$ to $\mathcal{A}^s(\Lambda; V^*)$ for all $s \in (0, s^*)$, see [CDD01]. This carries over to the routine Apply_R in that if $v \in \mathcal{A}^s(\Lambda; V)$ and z is the output of $\text{Apply}_R[v, \epsilon]$ for an $\epsilon > 0$, then

$$\#\text{supp } z \lesssim \|v\|_{\mathcal{A}^s(\Lambda; V)}^{1/s} \epsilon^{-1/s}, \quad (4.12)$$

$$\|z\|_{\mathcal{A}^s(\Lambda; V^*)} \lesssim \|v\|_{\mathcal{A}^s(\Lambda; V)} \quad (4.13)$$

with constants depending only on s and R . Moreover, (4.12) is an upper bound for the total number of applications of operators R_m in $\text{Apply}_R[v, \epsilon]$. This follows as in the standard scalar case, see *e.g.* [DSS09]. \dashv

We make additional assumptions on the routine RHS_f . If $f \in \mathcal{A}^s(\Lambda; V^*)$ and \tilde{f} is the output of $\text{RHS}_f[\epsilon]$ for an $\epsilon > 0$, then \tilde{f} should satisfy

$$\#\text{supp } \tilde{f} \lesssim \|f\|_{\mathcal{A}^s(\Lambda; V^*)}^{1/s} \epsilon^{-1/s}. \quad (4.14)$$

Note that if $u \in \mathcal{A}^s(\Lambda; V)$ and R is s^* -compressible with $s < s^*$, then also A is s^* -compressible, and therefore $\|f\|_{\mathcal{A}^s(\Lambda; V^*)} \lesssim \|u\|_{\mathcal{A}^s(\Lambda; V)}$.

Lemma 4.6. *Under the conditions of Theorem 4.2,*

$$\|\tilde{u}^{(k)}\|_{\mathcal{A}^s(\Lambda; V)} \leq C \|u\|_{\mathcal{A}^s(\Lambda; V)} \quad \forall k \in \mathbb{N}_0, \quad (4.15)$$

with

$$C = 1 + \frac{2^{1+s} \varrho (1 + \gamma) (1 + \omega)}{\tau (1 - \varrho^{1/s})^s (1 - \gamma) (1 - \omega)}, \quad (4.16)$$

$$\varrho = \sigma + \sqrt{1 - \vartheta^2 (1 - \gamma) (1 + \gamma)^{-1}} \text{ and } \tau = \sqrt{1 - \hat{\vartheta}^2 (1 + \gamma) (1 - \gamma)^{-1}}.$$

Proof. Let $k \in \mathbb{N}_0$. For any $N \geq \#\Xi^{(k)}$, $\|\tilde{u}^{(k)} - P_N(\tilde{u}^{(k)})\|_{\ell^2(\Lambda; V)} = 0$. For $N \leq \#\Xi^{(k)} - 1$,

$$\begin{aligned} \|\tilde{u}^{(k)} - P_N(\tilde{u}^{(k)})\|_{\ell^2(\Lambda; V)} &\leq \|\tilde{u}^{(k)} - \Pi_{\Xi_N} \tilde{u}^{(k)}\|_{\ell^2(\Lambda; V)} \\ &\leq \|u - \Pi_{\Xi_N} u\|_{\ell^2(\Lambda; V)} + 2 \|u - \tilde{u}^{(k)}\|_{\ell^2(\Lambda; V)}, \end{aligned}$$

where $\Xi_N := \text{supp } P_N(u)$, such that $\Pi_{\Xi_N} u = P_N(u)$ and

$$\|u - \Pi_{\Xi_N} u\|_{\ell^2(\Lambda; V)} \leq (N+1)^{-s} \|u\|_{\mathcal{A}^s(\Lambda; V)}.$$

Furthermore, Theorem 4.2 implies

$$\|\mathbf{u} - \tilde{\mathbf{u}}^{(k)}\|_{\ell^2(\Lambda; V)} \leq \frac{2^s \varrho(1 + \gamma)(1 + \omega)}{\tau(1 - \varrho^{1/s})^s(1 - \gamma)(1 - \omega)} (\#\Xi^{(k)})^{-s} \|\mathbf{u}\|_{\mathcal{A}^s(\Lambda; V)},$$

and $(\#\Xi^{(k)})^{-s} \leq (N + 1)^{-s}$ by definition of N . Consequently,

$$\|\tilde{\mathbf{u}}^{(k)}\|_{\mathcal{A}^s(\Lambda; V)} = \sup_{N \in \mathbb{N}_0} (N + 1)^{-s} \|\tilde{\mathbf{u}}^{(k)} - P_N(\tilde{\mathbf{u}}^{(k)})\|_{\ell^2(\Lambda; V)} \leq C \|\mathbf{u}\|_{\mathcal{A}^s(\Lambda; V)}$$

with C from (4.16). \square

Theorem 4.7. *Let the conditions of Theorem 4.2 be satisfied. If (4.14) and the assumptions of Proposition 4.4 hold for all $s \in (0, s^*)$, then for any $\epsilon > 0$ and any $s \in (0, s^*)$, the total number of applications of D , A_{VV} and D^{-1} in $\text{SolveGalerkin}_{\mathbf{A}, f}[\epsilon, \gamma, \chi, \vartheta, \omega, \sigma, 0, \beta]$ is bounded by $\|\mathbf{u}\|_{\mathcal{A}^s(\Lambda; V)}^{1/s} \epsilon^{-1/s}$ up to a constant factor depending only on the input arguments other than ϵ . The same bound holds for the total number of applications of R_m , $m \in \mathbb{N}$, up to an additional factor of $\max_{\mu \in \text{supp } \mathbf{u}_\epsilon} \#\text{supp } \mu$.*

Proof. Let $k \in \mathbb{N}_0$; we consider the k -th iteration of the loop in $\text{SolveGalerkin}_{\mathbf{A}, f}$. The routine $\text{Residual}_{\mathbf{A}, f}[\epsilon \sqrt{1 - \gamma}, \tilde{\mathbf{u}}^{(\Xi^{(k)})}, \delta_k, \chi, \omega, \beta]$ begins with $\#\Xi^{(k)}$ applications of D . Due to the geometric decrease in tolerances, the complexity of the loop in $\text{Residual}_{\mathbf{A}, f}$ is dominated by that of its last iteration. By Remark 4.5 and Lemma 4.6, the number of applications of D^{-1} and R_m is bounded by $\|\mathbf{u}\|_{\mathcal{A}^s(\Lambda; V)}^{1/s} \zeta_k^{-1/s}$, and $\zeta_k \gtrsim \bar{\delta}_k$.

Next, assuming the termination criterion of $\text{SolveGalerkin}_{\mathbf{A}, f}$ is not satisfied, the routine $\text{Galerkin}_{\mathbf{A}, f}[\Xi^{(k+1)}, \omega, \sigma \min(\delta_k, \bar{\delta}_k)]$ is called to iteratively approximate the Galerkin projection onto $\ell^2(\Xi^{(k+1)}; V)$. Since only a fixed relative error reduction is required, the number of iterations remains bounded. Therefore, the number of applications of D^{-1} and A_{VV} is bounded by $\#\Xi^{(k+1)}$ and the total number of applications of R_m , $m \in \mathbb{N}$, is bounded by $2\bar{\lambda}(\Xi^{(k+1)})\#\Xi^{(k+1)}$, where $\bar{\lambda}(\Xi^{(k+1)})$ denote the average length of indices in $\Xi^{(k+1)}$, see [Git11c, Proposition 3.5]. Since the sets $\Xi^{(k)}$ are nested, $\bar{\lambda}(\Xi^{(k+1)}) \leq \max_{\mu \in \text{supp } \mathbf{u}_\epsilon} \#\text{supp } \mu$. Furthermore, by Theorems 3.4 and 4.2, $\#\Xi^{(k+1)} \lesssim \|\mathbf{u}\|_{\mathcal{A}^s(\Lambda; V)}^{1/s} \bar{\delta}_{k+1}^{-1/s}$.

Let k be such that $\mathbf{u}_\epsilon = \tilde{\mathbf{u}}^{(k)}$. Due to the different termination criterion, the complexity of the last call of $\text{Residual}_{\mathbf{A}, f}$ can be estimated by $\|\mathbf{u}\|_{\mathcal{A}^s(\Lambda; V)}^{1/s} \zeta_k^{-1/s}$ with $\zeta_k \gtrsim \epsilon$. This bound obviously also holds for $\#\Xi^{(k)}$, and thus for the complexity of the final call of $\text{Galerkin}_{\mathbf{A}, f}$.

Combining all of the above estimates, the number of applications of D^{-1} , D , A_{VV} and R_m , $m \in \mathbb{N}$, in $\text{SolveGalerkin}_{\mathbf{A}, f}$ is bounded by

$$\|\mathbf{u}\|_{\mathcal{A}^s(\Lambda; V)}^{1/s} \left(\epsilon^{-1/s} + \sum_{j=0}^{k-1} \bar{\delta}_j^{-1/s} \right).$$

Furthermore, $\bar{\delta}_{k-1} \geq \epsilon$, and using $\delta_{k-1} \leq \varrho^{k-1-j} \bar{\delta}_j$ for $\varrho = \sigma + \sqrt{1 - \vartheta^2(1 - \gamma)(1 + \gamma)^{-1}} < 1$,

$$\sum_{j=0}^{k-2} \bar{\delta}_j^{-1/s} \leq \delta_{k-1}^{-1/s} \sum_{j=0}^{k-2} \varrho^{(k-1-j)/s} = \delta_{k-1}^{-1/s} \sum_{i=1}^{k-1} \varrho^{i/s} \leq \delta_{k-1}^{-1/s} \frac{\varrho^{1/s}}{1 - \varrho^{1/s}}.$$

The assertion follows since $\delta_{k-1} \geq \epsilon$. □

5 Computational Examples

5.1 Application to Isotropic Diffusion

We consider the isotropic diffusion equation on a bounded Lipschitz domain $G \subset \mathbb{R}^d$ with homogeneous Dirichlet boundary conditions. For any uniformly positive $a \in L^\infty(G)$ and any $f \in L^2(G)$, we have

$$\begin{aligned} -\nabla \cdot (a(x)\nabla u(x)) &= f(x), \quad x \in G, \\ u(x) &= 0, \quad x \in \partial G. \end{aligned} \quad (5.1)$$

We view f as fixed, but allow a to vary, giving rise to a parametric operator

$$A_0(a): H_0^1(G) \rightarrow H^{-1}(G), \quad v \mapsto -\nabla \cdot (a\nabla v), \quad (5.2)$$

which depends continuously on $a \in L^\infty(G)$.

We model the coefficient a as a bounded $L^\infty(G)$ -valued random field, which we expand as a series

$$a(y, x) := \bar{a}(x) + \sum_{m=1}^{\infty} y_m a_m(x). \quad (5.3)$$

Since a is bounded, a_m can be scaled such that $y_m \in [-1, 1]$ for all $m \in \mathbb{N}$. Therefore, a depends on a parameter $y = (y_m)_{m=1}^{\infty}$ in $\Gamma = [-1, 1]^\infty$.

We define the parametric operator $A(y) := A_0(a(y))$ for $y \in \Gamma$. Due to the linearity of A_0 ,

$$A(y) = D + R(y), \quad R(y) := \sum_{m \in \mathcal{M}} y_m R_m \quad \forall y \in \Gamma \quad (5.4)$$

with convergence in $\mathcal{L}(H_0^1(G), H^{-1}(G))$, for

$$\begin{aligned} D &:= A_0(\bar{a}): H_0^1(G) \rightarrow H^{-1}(G), \quad v \mapsto -\nabla \cdot (\bar{a}\nabla v), \\ R_m &:= A_0(a_m): H_0^1(G) \rightarrow H^{-1}(G), \quad v \mapsto -\nabla \cdot (a_m\nabla v), \quad m \in \mathcal{M}. \end{aligned}$$

To ensure bounded invertibility of D , we assume there is a constant $\delta > 0$ such that

$$\operatorname{ess\,inf}_{x \in G} \bar{a}(x) \geq \delta^{-1}. \quad (5.5)$$

We refer *e.g.* to [Git11c, Git11a, SG11] for further details.

5.2 A Posteriori Error Estimation

In $\text{SolveGalerkin}_{A,f}$, a generic solver Solve_D is used to approximate $D^{-1}g_\mu$ to any desired accuracy, where g_μ has the form

$$g_\mu = f_\mu - \sum_{i=1}^k \kappa_i R_{m_i} w_i, \quad (5.6)$$

with $w_i \in V = H_0^1(G)$ equal to some coefficients \tilde{w}_i of the previous approximate solution. If $D^{-1}g_\mu$ is approximated by the finite element method, an a posteriori error estimator is required to determine whether or not a given approximation attains the desired accuracy. Due to the unusual structure of g_μ , standard error estimators cannot be applied directly. We derive a reliable residual-based estimator, following the standard argument from [MNS00, AO00, Ver96].

Let \mathfrak{T} be a regular mesh of G , and let W_μ be a finite element space of continuous, piecewise smooth shape functions on \mathfrak{T} which contains at least the piecewise linear functions.

We will denote the set of elements of \mathfrak{T} by \mathfrak{T} and the set of faces of \mathfrak{T} by \mathfrak{F} . The set \mathfrak{F} can be decomposed into interior faces $\mathfrak{F} \cap G$ and boundary faces $\mathfrak{F} \cap \partial G$. For any $T \in \mathfrak{T}$, let h_T be the diameter of T , and similarly, define h_F as the diameter of F for any $F \in \mathfrak{F}$. Furthermore, for any $T \in \mathfrak{T}$, let $\tilde{\omega}_T \subset G$ consist of all elements of \mathfrak{T} sharing at least a vertex with T . Analogously, let $\tilde{\omega}_F \subset G$ consist of all elements of \mathfrak{T} sharing at least a vertex with the face $F \in \mathfrak{F}$. Note that each element $T \in \mathfrak{T}$ belongs to only a bounded number of domains $\tilde{\omega}_{T'}$ or $\tilde{\omega}_F$.

By the above assumptions, there is a Clément interpolant for W_μ , i.e. a continuous projection $\mathfrak{I}_\mu: H_0^1(G) \rightarrow W_\mu$ such that for all $v \in H_0^1(G)$,

$$\|v - \mathfrak{I}_\mu v\|_{L^2(T)} \leq c_1 h_T |v|_{H^1(\tilde{\omega}_T)} \quad \forall T \in \mathfrak{T} \quad (5.7)$$

and

$$\|v - \mathfrak{I}_\mu v\|_{L^2(F)} \leq c_2 h_F^{1/2} |v|_{H^1(\tilde{\omega}_F)} \quad \forall F \in \mathfrak{F} \quad (5.8)$$

with constants c_1 and c_2 depending only on the shape regularity of \mathfrak{T} , see e.g. [BS02].

Let each of the functions w_i from (5.6) itself be an element of a finite element space W_i of piecewise smooth functions on a mesh \mathfrak{T}_i , which may differ from \mathfrak{T} . We assume that these meshes are compatible in the sense that for any $T \in \mathfrak{T}$ and $T_i \in \mathfrak{T}_i$, the intersection $T \cap T_i$ is either empty, equal to T , or equal to T_i .

Standard error estimators run into problems on faces of \mathfrak{T}_i that are not in the skeleton of \mathfrak{T} , since g_μ is singular on these faces. For all i , let \bar{w}_i be an approximation of w_i that is piecewise smooth on \mathfrak{T} . Replacing g_μ by

$$\bar{g}_\mu := f_\mu - \sum_{i=1}^k \kappa_i R_{m_i} \bar{w}_i \quad (5.9)$$

induces an error

$$\|D^{-1}g_\mu - D^{-1}\bar{g}_\mu\|_V \leq \sum_{i=1}^k |\kappa_i| \left\| \frac{a_{m_i}}{\bar{a}} \right\|_{L^\infty(G)} \|w_i - \bar{w}_i\|_V =: \text{EST}_\mu^P, \quad (5.10)$$

since

$$\sup_{\|v\|_V=1} \left| \int_G a_m \nabla w \cdot \nabla v \, dx \right| \leq \left\| \frac{a_m}{\bar{a}} \right\|_{L^\infty(G)} \sup_{\|v\|_V=1} \int_G |\bar{a} \nabla w \cdot \nabla v| \, dx = \left\| \frac{a_m}{\bar{a}} \right\|_{L^\infty(G)} \|w\|_V$$

for all $m \in \mathbb{N}$ and all $w \in H_0^1(G)$.

Let $\bar{u}_\mu \in W_\mu$ be the Galerkin projection of $D^{-1}\bar{g}_\mu$, i.e.

$$\int_G \bar{a} \nabla \bar{u}_\mu \cdot \nabla v \, dx = \int_G f_\mu v \, dx - \sum_{i=1}^k \kappa_i \int_G a_{m_i} \nabla \bar{w}_i \cdot \nabla v \, dx \quad \forall v \in W_\mu. \quad (5.11)$$

Abbreviating

$$\sigma_\mu := \bar{a} \nabla \bar{u}_\mu + \sum_{i=1}^k \kappa_i a_{m_i} \nabla \bar{w}_i, \quad (5.12)$$

the residual of \bar{u}_μ is the functional

$$r_\mu(\bar{u}_\mu; v) = \int_G \bar{g}_\mu - \bar{a} \nabla \bar{u}_\mu \cdot \nabla v \, dx = \int_G f_\mu - \sigma_\mu \cdot \nabla v \, dx, \quad v \in H_0^1(G). \quad (5.13)$$

By Galerkin orthogonality, $r_\mu(\bar{u}_\mu; v) = 0$ for all $v \in W_\mu$. Furthermore, due to the Riesz isomorphism,

$$\|D^{-1}\bar{g}_\mu - \bar{u}_\mu\|_V = \sup_{v \in H_0^1(G) \setminus \{0\}} \frac{|r_\mu(\bar{u}_\mu; v)|}{\|v\|_D} \leq \sqrt{\delta} \sup_{v \in H_0^1(G) \setminus \{0\}} \frac{|r_\mu(\bar{u}_\mu; v)|}{\|v\|_{H^1(G)}}, \quad (5.14)$$

with δ from (5.5).

For all $T \in \mathfrak{T}$, let

$$R_{\mu,T}(\bar{u}_\mu) := h_T \|f_\mu + \nabla \cdot \sigma_\mu\|_{L^2(T)}, \quad (5.15)$$

where the dependence on \bar{u}_μ is implicit in σ_μ . Note that $\nabla \cdot \sigma_\mu$ is given by

$$\nabla \cdot \sigma_\mu = \nabla \bar{a} \cdot \nabla \bar{u}_\mu + \bar{a} \Delta \bar{u}_\mu + \sum_{i=1}^k \kappa_i (\nabla a_{m_i} \cdot \nabla \bar{w}_i + a_{m_i} \Delta \bar{w}_i). \quad (5.16)$$

Also, let

$$R_{\mu,F}(\bar{u}_\mu) := h_F^{1/2} \|[\![\sigma_\mu]\!] \|_{L^2(F)}, \quad (5.17)$$

where $[\![\cdot]\!]$ is the normal jump over the face $F \in \mathfrak{F} \cap G$, i.e. if $F = T_1 \cap T_2$, and n_1 and n_2 are the respective exterior normal vectors, then

$$[\![\sigma_\mu]\!] := \sigma_\mu|_{T_1} \cdot n_1 + \sigma_\mu|_{T_2} \cdot n_2, \quad (5.18)$$

and $[\![\sigma_\mu]\!] := \sigma_\mu \cdot n_G$ if $F \in \mathfrak{F} \cap \partial G$ for the exterior unit normal n_G of G . These terms combine to

$$\text{EST}_\mu^R(\bar{u}_\mu) := \left(\sum_{T \in \mathfrak{T}} R_{\mu,T}(\bar{u}_\mu)^2 + \sum_{F \in \mathfrak{F}} R_{\mu,F}(\bar{u}_\mu)^2 \right)^{1/2}. \quad (5.19)$$

Note that if $d = 1$, then $h_F = 0$ for all $F \in \mathfrak{F}$, and $R_{\mu,F}(\bar{u}_\mu) = 0$. In this case, the Clément interpolation operator \mathfrak{I}_μ is simply the nodal interpolant, and $\bar{\omega}_T$ can be replaced by T in (5.7).

Theorem 5.1. For all $v \in H_0^1(G)$,

$$|r_\mu(\bar{u}_\mu; v)| \leq C \text{EST}_\mu^R(\bar{u}_\mu) |v|_{H^1(G)} \quad (5.20)$$

with a constant C depending only on the shape regularity of \mathfrak{T} .

Proof. Let $v \in H_0^1(G)$. Since $\mathfrak{S}_\mu v \in W_\mu$, by Galerkin orthogonality

$$r_\mu(\bar{u}_\mu; v) = r_\mu(\bar{u}_\mu; v - \mathfrak{S}_\mu v) .$$

We abbreviate $w := v - \mathfrak{S}_\mu v$, and denote by n_T the exterior unit normal of $T \in \mathfrak{T}$. Using (5.7), (5.8), integration by parts and the Cauchy–Schwarz inequality,

$$\begin{aligned} r_\mu(\bar{u}_\mu; v - \mathfrak{S}_\mu v) &= \sum_{T \in \mathfrak{T}} \int_T f_\mu w - \sigma_\mu \cdot \nabla w \, dx \\ &= \sum_{T \in \mathfrak{T}} \left[\int_T (f_\mu + \nabla \cdot \sigma_\mu) w \, dx - \sum_{F \in \mathfrak{F} \cap \partial T} \int_F \sigma_\mu \cdot n_T w \, dS \right] \\ &= \sum_{T \in \mathfrak{T}} \int_T (f_\mu + \nabla \cdot \sigma_\mu) w \, dx - \sum_{F \in \mathfrak{F}} \int_F \llbracket \sigma_\mu \rrbracket w \, dS \\ &\leq \sum_{T \in \mathfrak{T}} \|f_\mu + \nabla \cdot \sigma_\mu\|_{L^2(T)} \|w\|_{L^2(T)} + \sum_{F \in \mathfrak{F}} \|\llbracket \sigma_\mu \rrbracket\|_{L^2(F)} \|w\|_{L^2(F)} \\ &\leq c_1 \sum_{T \in \mathfrak{T}} h_T \|f_\mu + \nabla \cdot \sigma_\mu\|_{L^2(T)} |v|_{H^1(\bar{\omega}_T)} + c_2 \sum_{F \in \mathfrak{F}} h_F^{1/2} \|\llbracket \sigma_\mu \rrbracket\|_{L^2(F)} |v|_{H^1(\bar{\omega}_F)} \\ &\leq (c_1 + c_2) \sum_{T \in \mathfrak{T}} \left[h_T \|f_\mu + \nabla \cdot \sigma_\mu\|_{L^2(T)} + \sum_{F \in \mathfrak{F} \cap \partial T} h_F^{1/2} \|\llbracket \sigma_\mu \rrbracket\|_{L^2(F)} \right] |v|_{H^1(\bar{\omega}_T)} \\ &\leq C_0 \left(\sum_{T \in \mathfrak{T}} \left[h_T \|f_\mu + \nabla \cdot \sigma_\mu\|_{L^2(T)} + \sum_{F \in \mathfrak{F} \cap \partial T} h_F^{1/2} \|\llbracket \sigma_\mu \rrbracket\|_{L^2(F)} \right]^2 \right)^{1/2} |v|_{H^1(G)} \\ &\leq C \text{EST}_\mu^R(\bar{u}_\mu) |v|_{H^1(G)} . \end{aligned}$$

This shows (5.20), replacing v by $-v$ if necessary. \square

Corollary 5.2. The Galerkin projection \bar{u}_μ from (5.11) satisfies

$$\|D^{-1}g_\mu - \bar{u}_\mu\|_V \leq \text{EST}_\mu^P + \sqrt{\delta} C \text{EST}_\mu^R(\bar{u}_\mu) \quad (5.21)$$

for δ from (5.5) and C from Theorem 5.1.

Proof. The assertion follows by triangle inequality using (5.10), (5.14) and (5.20). \square

5.3 Numerical Computations

We consider as a model problem the diffusion equation (5.1) on the one dimensional domain $G = (0, 1)$. For two parameters k and γ , the diffusion coefficient has the form

$$a(y, x) = 1 + \frac{1}{c} \sum_{m=1}^{\infty} y_m \frac{1}{m^k} \sin(m\pi x), \quad x \in (0, 1), \quad y \in \Gamma = [-1, 1]^{\infty}, \quad (5.22)$$

where c is chosen as

$$c = \gamma \sum_{m=1}^{\infty} \frac{1}{m^k}, \quad (5.23)$$

such that $|a(y, x) - 1|$ is always less than γ . For the distribution of $y \in \Gamma$, we consider the countable product of uniform distributions on $[-1, 1]$; the corresponding family of orthonormal polynomials is the Legendre polynomial basis.

In all of the following computations, the parameters are $k = 2$ and $\gamma = 1/2$. A few realizations of $a(y)$ and the resulting solutions $u(y)$ of (5.1) are plotted in Figure 1.

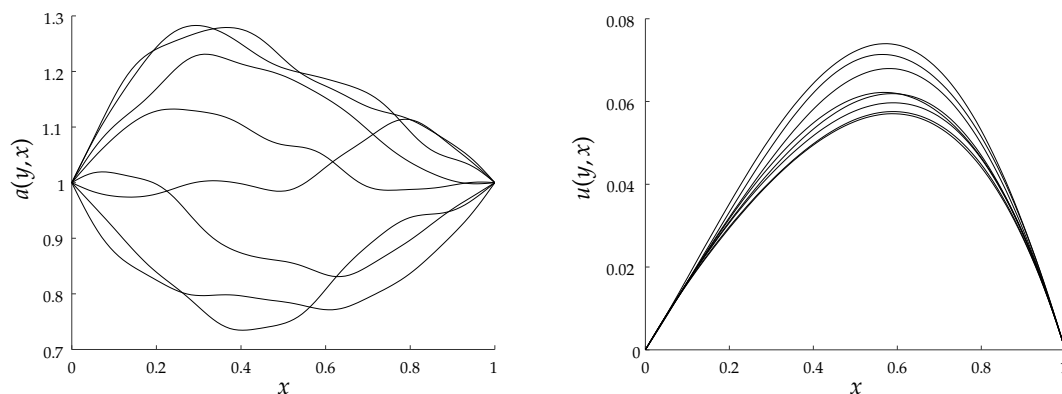


Figure 1: Realizations of $a(y, x)$ (left) and $u(y, x)$ (right).

The parameters of `SolveGalerkinA,f` are set to $\chi = 1/8$, $\vartheta = 0.57$, $\omega = 1/4$, $\sigma = 0.01114$, $\alpha = 1/20$ and $\beta = 0$. These values do not satisfy the assumptions of Theorem 4.2; however, the method executes substantially faster than with parameters for which the theorem applies. All computations were performed in Matlab on a workstation with an AMD Athlon™ 64 X2 5200+ processor and 4GB of memory.

We consider a multilevel discretization in which the a posteriori error estimator from Section 5.2 is used to determine an appropriate discretization level for each coefficient. A discretization level j_μ , which represents linear finite elements on a uniform mesh with 2^{j_μ} cells, is assigned to each index μ with the goal of equidistributing the estimated error among all coefficients.

In Figure 2, on the left, the errors are plotted against the number of degrees of freedom, which refers to the total number of basis functions used in the discretization, *i.e.* the sum of $2^{j_\mu} - 1$ over all μ . On the right, we plot the errors against an estimate of the

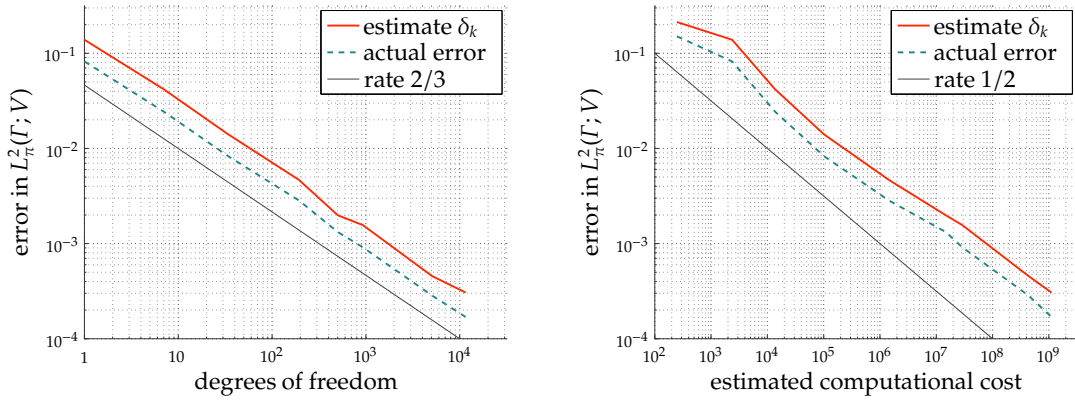


Figure 2: Convergence of $\text{SolveGalerkin}_{A,f}$.

computational cost. This estimate takes scalar products, matrix-vector multiplications and linear solves into account. The total number of each of these operations on each discretization level is tabulated during the computation, weighted by the number of degrees of freedom on the discretization level, and summed over all levels. The estimate is equal to seven times the resulting sum for linear solves, plus three times the value for matrix-vector multiplications, plus the sum for scalar products. These weights were determined empirically by timing the operations for tridiagonal sparse matrices in Matlab.

The errors were computed by comparison with a reference solution, which has an error of approximately $5 \cdot 10^{-5}$. The plots show that the error bounds δ_k are good approximations of the error, and only overestimate it by a small factor.

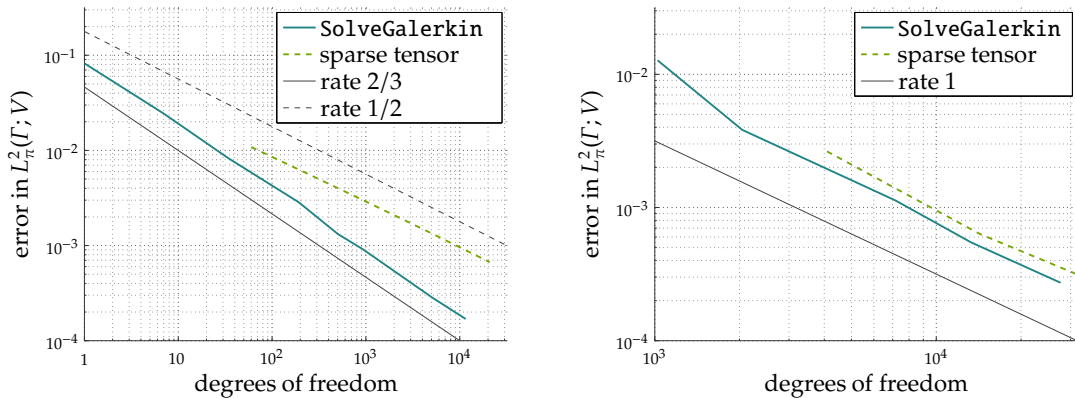


Figure 3: Comparison of $\text{SolveGalerkin}_{A,f}$ and the sparse tensor construction, for a multilevel discretization (left) and with a fixed finite element mesh (right).

We compare the discretizations generated adaptively by $\text{SolveGalerkin}_{A,f}$ with the heuristic a priori adapted sparse tensor product construction from [BAS10]. Using the notation of [SG11, Section 4], we set $\gamma = 2$ and $\eta_m = 1/(r_m + \sqrt{1 + r_m^2})$ for $r_m = cm^2/2$ and

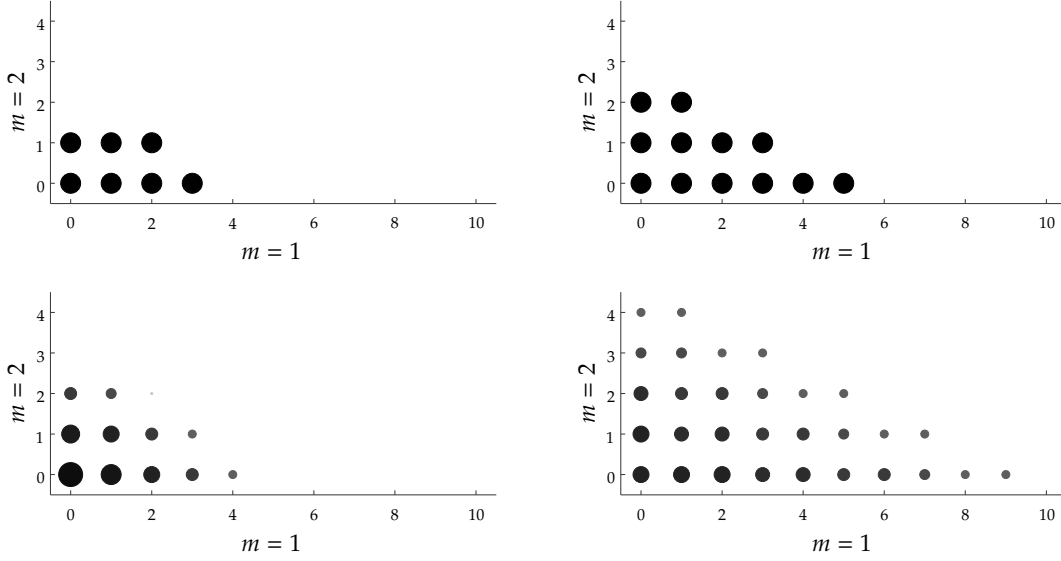


Figure 4: Slices of index sets generated by $\text{SolveGalerkin}_{A,f}$ (left) and [BAS10] (right) with single level discretization (top) and multilevel discretization (bottom). All sets correspond to the right-most points in Figure 3. Active indices with support in $\{1,2\}$ are plotted; the level of the finite element discretization is proportional to the radius of the circle.

c from (5.23). These values are similar to those used in the computational examples of [BAS10]. The coarsest spatial discretization used in the sparse tensor product contains 16 elements.

In order to isolate the stochastic discretization, we also consider a fixed spatial discretization, using linear finite elements on a uniform mesh of $(0,1)$ with 1024 elements to approximate all coefficients. We refer to these simpler versions of the numerical methods as single level discretizations.

The single level versions of $\text{SolveGalerkin}_{A,f}$ and the sparse tensor method construct discretizations of equal quality, with only a slight advantage for the adaptive algorithm. However, with a multilevel discretization, $\text{SolveGalerkin}_{A,f}$ converges faster than the sparse tensor method, with respect to the number of degrees of freedom. At least in this example, the adaptively constructed discretizations are more efficient than sparse tensor products.

As index sets $\mathcal{E} \subset \Lambda$ are infinite dimensional in the sense that they can contain indices of arbitrary length, they are difficult to visualize in only two dimensions. In Figure 4, we plot two dimensional slices of sets generated by $\text{SolveGalerkin}_{A,f}$ and the sparse tensor construction from [BAS10]. We consider only those indices which are zero in all dimensions after the second, and plot their values in the first two dimensions. The upper plots depict index sets generated using single level discretizations; dots refer to active indices. The lower plots illustrate the discretizations generated with multilevel finite element discretizations. The radii of the circles are proportional to the discretization

level.

The bottom two plots in Figure 4 illustrate differences between the discretizations generated by $\text{SolveGalerkin}_{A,f}$ and the sparse tensor construction. The former has many fewer active indices, but higher discretization levels for some of these. For example, the coefficient of the constant polynomial is approximated on meshes with 4096 and 256 elements, respectively. Also, while the sets constructed by sparse tensorization appear triangular in this figure, the adaptively generated index sets are somewhat more convex. All of the sets are anisotropic in the sense that the first dimension is discretized more finely than the second.

We use the convergence curves in Figures 2 and 3 to empirically determine convergence rates of $\text{SolveGalerkin}_{A,f}$. The convergence rate with respect to the total number of degrees of freedom is $2/3$, which is faster than the approximation of $1/2$ rate shown in [CDS10b, CDS10a]. It also compares favorably to the sparse tensor construction, which converges at a rate of $1/2$. However, when considering convergence with respect to the computational cost, the rate of $\text{SolveGalerkin}_{A,f}$ reduces to $1/2$ also. We suspect that this is due to the approximation of the residual, which is performed on a larger set of active indices than the subsequent approximation of the Galerkin projection.

The solvers with fixed finite element meshes simulate semi-discrete methods with no spatial discretization. In this setting, [CDS10b, CDS10a] show an approximation rate of $3/2$, whereas we observe a rate of 1 for both $\text{SolveGalerkin}_{A,f}$ and sparse tensorization. In principle, it is possible that $\text{SolveGalerkin}_{A,f}$ does not converge with the optimal rate in this example, since the parameters used in the computations do not satisfy the assumptions of Theorem 4.2. Alternatively, due to large constants in the approximation estimates, the asymptotic rate may not be perceivable for computationally accessible tolerances.

References

- [AO00] Mark Ainsworth and J. Tinsley Oden. *A posteriori error estimation in finite element analysis*. Pure and Applied Mathematics (New York). Wiley-Interscience [John Wiley & Sons], New York, 2000.
- [Bar05] A. Barinka. *Fast Evaluation Tools for Adaptive Wavelet Schemes*. PhD thesis, RWTH Aachen, March 2005.
- [BAS10] Marcel Bieri, Roman Andreev, and Christoph Schwab. Sparse tensor discretization of elliptic SPDEs. *SIAM J. Sci. Comput.*, 31(6):4281–4304, 2009/10.
- [Bau02] Heinz Bauer. *Wahrscheinlichkeitstheorie*. de Gruyter Lehrbuch. [de Gruyter Textbook]. Walter de Gruyter & Co., Berlin, fifth edition, 2002.
- [BDD04] Peter Binev, Wolfgang Dahmen, and Ron DeVore. Adaptive finite element methods with convergence rates. *Numer. Math.*, 97(2):219–268, 2004.

- [BS02] Susanne C. Brenner and L. Ridgway Scott. *The mathematical theory of finite element methods*, volume 15 of *Texts in Applied Mathematics*. Springer-Verlag, New York, second edition, 2002.
- [BS09] Marcel Bieri and Christoph Schwab. Sparse high order FEM for elliptic sPDEs. *Comput. Methods Appl. Mech. Engrg.*, 198(37-40):1149–1170, 2009.
- [BTZ04] Ivo Babuška, Raúl Tempone, and Georgios E. Zouraris. Galerkin finite element approximations of stochastic elliptic partial differential equations. *SIAM J. Numer. Anal.*, 42(2):800–825 (electronic), 2004.
- [CDD01] Albert Cohen, Wolfgang Dahmen, and Ronald DeVore. Adaptive wavelet methods for elliptic operator equations: convergence rates. *Math. Comp.*, 70(233):27–75 (electronic), 2001.
- [CDD02] A. Cohen, W. Dahmen, and R. DeVore. Adaptive wavelet methods. II. Beyond the elliptic case. *Found. Comput. Math.*, 2(3):203–245, 2002.
- [CDS10a] A. Cohen, R. DeVore, and Ch. Schwab. Analytic regularity and polynomial approximation of parametric stochastic elliptic PDEs. Technical Report 2010-3, Seminar for Applied Mathematics, ETH Zürich, 2010. In review.
- [CDS10b] Albert Cohen, Ronald DeVore, and Christoph Schwab. Convergence rates of best N -term Galerkin approximations for a class of elliptic sPDEs. *Found. Comput. Math.*, 10(6):615–646, 2010.
- [DBO01] Manas K. Deb, Ivo M. Babuška, and J. Tinsley Oden. Solution of stochastic partial differential equations using Galerkin finite element techniques. *Comput. Methods Appl. Mech. Engrg.*, 190(48):6359–6372, 2001.
- [DeV98] Ronald A. DeVore. Nonlinear approximation. In *Acta numerica, 1998*, volume 7 of *Acta Numer.*, pages 51–150. Cambridge Univ. Press, Cambridge, 1998.
- [DFR07] Stephan Dahlke, Massimo Fornasier, and Thorsten Raasch. Adaptive frame methods for elliptic operator equations. *Adv. Comput. Math.*, 27(1):27–63, 2007.
- [Dör96] Willy Dörfler. A convergent adaptive algorithm for Poisson’s equation. *SIAM J. Numer. Anal.*, 33(3):1106–1124, 1996.
- [DRW⁺07] Stephan Dahlke, Thorsten Raasch, Manuel Werner, Massimo Fornasier, and Rob Stevenson. Adaptive frame methods for elliptic operator equations: the steepest descent approach. *IMA J. Numer. Anal.*, 27(4):717–740, 2007.
- [DSS09] Tammo Jan Dijkema, Christoph Schwab, and Rob Stevenson. An adaptive wavelet method for solving high-dimensional elliptic PDEs. *Constr. Approx.*, 30(3):423–455, 2009.

- [FST05] Philipp Frauenfelder, Christoph Schwab, and Radu Alexandru Todor. Finite elements for elliptic problems with stochastic coefficients. *Comput. Methods Appl. Mech. Engrg.*, 194(2-5):205–228, 2005.
- [GHS07] Tsogtgerel Gantumur, Helmut Harbrecht, and Rob Stevenson. An optimal adaptive wavelet method without coarsening of the iterands. *Math. Comp.*, 76(258):615–629 (electronic), 2007.
- [Git11a] Claude Jeffrey Gittelsohn. *Adaptive Galerkin Methods for Parametric and Stochastic Operator Equations*. PhD thesis, ETH Zürich, 2011. ETH Dissertation No. 19533.
- [Git11b] Claude Jeffrey Gittelsohn. Adaptive stochastic Galerkin methods: Beyond the elliptic case. Technical Report 2011-12, Seminar for Applied Mathematics, ETH Zürich, 2011.
- [Git11c] Claude Jeffrey Gittelsohn. Stochastic Galerkin approximation of operator equations with infinite dimensional noise. Technical Report 2011-10, Seminar for Applied Mathematics, ETH Zürich, 2011.
- [Met02] A. Metselaar. *Handling Wavelet Expansions in Numerical Methods*. PhD thesis, University of Twente, 2002.
- [MK05] Hermann G. Matthies and Andreas Keese. Galerkin methods for linear and nonlinear elliptic stochastic partial differential equations. *Comput. Methods Appl. Mech. Engrg.*, 194(12-16):1295–1331, 2005.
- [MNS00] Pedro Morin, Ricardo H. Nochetto, and Kunibert G. Siebert. Data oscillation and convergence of adaptive FEM. *SIAM J. Numer. Anal.*, 38(2):466–488 (electronic), 2000.
- [SG11] Christoph Schwab and Claude Jeffrey Gittelsohn. Sparse tensor discretization of high-dimensional parametric and stochastic PDEs. *Acta Numerica*, 2011. To appear.
- [Ste03] Rob Stevenson. Adaptive solution of operator equations using wavelet frames. *SIAM J. Numer. Anal.*, 41(3):1074–1100 (electronic), 2003.
- [TS07] Radu Alexandru Todor and Christoph Schwab. Convergence rates for sparse chaos approximations of elliptic problems with stochastic coefficients. *IMA J. Numer. Anal.*, 27(2):232–261, 2007.
- [Ver96] R. Verfürth. *A Review of a Posteriori Error Estimation and Adaptive Mesh-Refinement Techniques*. Teubner Verlag and J. Wiley, Stuttgart, 1996.
- [WK05] Xiaoliang Wan and George Em Karniadakis. An adaptive multi-element generalized polynomial chaos method for stochastic differential equations. *J. Comput. Phys.*, 209(2):617–642, 2005.

- [WK06] Xiaoliang Wan and George Em Karniadakis. Multi-element generalized polynomial chaos for arbitrary probability measures. *SIAM J. Sci. Comput.*, 28(3):901–928 (electronic), 2006.
- [XK02] Dongbin Xiu and George Em Karniadakis. The Wiener-Askey polynomial chaos for stochastic differential equations. *SIAM J. Sci. Comput.*, 24(2):619–644 (electronic), 2002.

Research Reports

No.	Authors/Title
11-11	<i>C.J. Gittelsohn</i> An adaptive stochastic Galerkin method
11-10	<i>C.J. Gittelsohn</i> Stochastic Galerkin approximation of operator equations with infinite dimensional noise
11-09	<i>R. Hiptmair, A. Moiola and I. Perugia</i> Error analysis of Trefftz-discontinuous Galerkin methods for the time-harmonic Maxwell equations
11-08	<i>W. Dahmen, C. Huang, Ch. Schwab and G. Welper</i> Adaptive Petrov-Galerkin methods for first order transport equations
11-07	<i>V.H. Hoang and Ch. Schwab</i> Analytic regularity and polynomial approximation of stochastic, parametric elliptic multiscale PDEs
11-06	<i>G.M. Coclite, K.H. Karlsen, S. Mishra and N.H. Risebro</i> A hyperbolic-elliptic model of two-phase flow in porous media - Existence of entropy solutions
11-05	<i>U.S. Fjordholm, S. Mishra and E. Tadmor</i> Entropy stable ENO scheme
11-04	<i>M. Ganesh, S.C. Hawkins and R. Hiptmair</i> Convergence analysis with parameter estimates for a reduced basis acoustic scattering T-matrix method
11-03	<i>O. Reichmann</i> Optimal space-time adaptive wavelet methods for degenerate parabolic PDEs
11-02	<i>S. Mishra, Ch. Schwab and J. Šukys</i> Multi-level Monte Carlo finite volume methods for nonlinear systems of conservation laws in multi-dimensions
11-01	<i>V. Wheatley, R. Jeltsch and H. Kumar</i> Spectral performance of RKDG methods
10-49	<i>R. Jeltsch and H. Kumar</i> Three dimensional plasma arc simulation using resistive MHD
10-48	<i>M. Swärd and S. Mishra</i> Entropy stable schemes for initial-boundary-value conservation laws