



Norwegian University of  
Science and Technology

# Reduced Basis Methods for Partial Differential Equations

Evaluation of multiple non-compliant flux-type output functionals  
for a non-affine electrostatics problem

**Jens Lohne Eftang**

Master of Science in Physics and Mathematics

Submission date: June 2008

Supervisor: Einar Rønquist, MATH



# Problem Description

The purpose of this work is to study reduced basis approximation methods for parametrised partial differential equations. Focus should be put on numerical accuracy and computational efficiency, and the report should include one or more illustrative examples.

Assignment given: 15. January 2008  
Supervisor: Einar Rønquist, MATH



# Preface

This thesis concludes my five-year study of physics and applied mathematics at the Norwegian University of Science and Technology (NTNU) in Trondheim, Norway. It constitutes my work during the last semester – an equivalent of 30 ECTS credits – and sorts under the Department of Mathematical Sciences, Faculty of Information Technology, Mathematics and Electrical Engineering with code TMA4900.

The implementations of the various numerical methods are all carried out in Matlab – a programming environment very convenient for rapid implementation and demonstration of ideas – as my interests have been in the methodology and not in optimising the code for runtime speed.

I would like to thank my supervisor, Professor Einar M. Rønquist, at the Department of Mathematical Sciences, NTNU. His good advice, many ideas and encouraging words have meant a lot during my work on this project.

Trondheim, 12th June 2008  
*Jens Lohne Eftang*



# Summary

A method for rapid evaluation of flux-type outputs of interest from solutions to partial differential equations (PDEs) is presented within the reduced basis framework for linear, elliptic PDEs. The central point is a Neumann-Dirichlet equivalence that allows for evaluation of the output through the bilinear form of the weak formulation of the PDE.

Through a comprehensive example related to electrostatics, we consider multiple outputs, *a posteriori* error estimators and empirical interpolation treatment of the non-affine terms in the bilinear form. Together with the considered Neumann-Dirichlet equivalence, these methods allow for efficient and accurate numerical evaluation of a relationship  $\boldsymbol{\mu} \rightarrow s(\boldsymbol{\mu})$ , where  $\boldsymbol{\mu}$  is a parameter vector that determines the geometry of the physical domain and  $s(\boldsymbol{\mu})$  is the corresponding flux-type output *matrix* of interest.

As a practical application, we lastly employ the rapid evaluation of  $\boldsymbol{\mu} \rightarrow s(\boldsymbol{\mu})$  in solving an inverse (parameter-estimation) problem.





# Contents

<b>Preface</b>	<b>i</b>
<b>Summary</b>	<b>iii</b>
<b>Contents</b>	<b>v</b>
<b>1 Introduction</b>	<b>1</b>
1.1 A brief intro to the reduced basis method . . . . .	1
1.2 Scope and overview . . . . .	2
1.3 A few remarks on notation . . . . .	3
<b>2 Preliminaries</b>	<b>5</b>
2.1 Weak form of the Poisson problem . . . . .	5
2.2 Norms and inner-products . . . . .	6
2.3 Gauss-Lobatto-Legendre Quadrature . . . . .	7
2.4 Existence and uniqueness of a weak solution . . . . .	7
2.5 Evaluation of flux-type output functionals . . . . .	8
<b>3 The Spectral Element (SE) method</b>	<b>15</b>
3.1 Spectral element discretisation . . . . .	15
3.2 Basis and algebraic formulation . . . . .	17
3.3 Implementation notes and operation count . . . . .	17
3.4 <i>A priori</i> error estimates . . . . .	18
3.4.1 Error in the field variable . . . . .	18
3.4.2 Error in the output of interest . . . . .	20
3.5 Numerical examples: Two model problems . . . . .	21
3.5.1 A problem with known analytic solution . . . . .	22
3.5.2 A problem with known singularity solution . . . . .	23
<b>4 Reduced Basis (RB) Approximation</b>	<b>27</b>
4.1 Formulation . . . . .	27
4.1.1 Parametric weak form . . . . .	27

4.1.2	Norms and inner-products . . . . .	28
4.1.3	Truth approximation . . . . .	29
4.1.4	Discrete formulation . . . . .	29
4.1.5	Algebraic formulation . . . . .	31
4.2	Offline-online procedure for affine problems . . . . .	33
4.3	Snapshot sampling, a greedy algorithm . . . . .	35
4.4	<i>A posteriori</i> error estimation . . . . .	36
4.4.1	Energy norm error bound . . . . .	36
4.4.2	Output error bounds . . . . .	37
<b>5</b>	<b>Empirical Interpolation (EI)</b>	<b>39</b>
5.1	Motivation . . . . .	39
5.2	Interpolation algorithm . . . . .	41
5.2.1	A remark on practical implementation . . . . .	43
5.3	Error estimation . . . . .	44
5.4	Numerical examples . . . . .	45
5.4.1	A one-dimensional example . . . . .	45
<b>6</b>	<b>A worked example: Electrostatics</b>	<b>49</b>
6.1	Physical principles . . . . .	50
6.1.1	The governing equation and boundary conditions . . . . .	50
6.1.2	Symmetry considerations . . . . .	51
6.2	RB treatment of the forward problem . . . . .	53
6.2.1	Parametric weak form . . . . .	53
6.2.2	Spectral element truth approximation . . . . .	54
6.2.3	RB formulation . . . . .	58
6.2.4	RB formulation with Empirical Interpolation (RB-EI) . . . . .	59
6.2.5	Remarks on RB and RB-EI output evaluation . . . . .	61
6.2.6	Remarks on inhomogeneous Dirichlet conditions . . . . .	64
6.2.7	<i>A posteriori</i> error estimation . . . . .	65
6.3	An inverse problem . . . . .	67
6.4	Numerical results . . . . .	69
6.4.1	Spectral element truth approximation . . . . .	69
6.4.2	Reduced basis approximation . . . . .	72
6.4.3	Parameter estimation . . . . .	79
<b>7</b>	<b>Conclusions</b>	<b>81</b>
	<b>Bibliography</b>	<b>83</b>

# Chapter 1

## Introduction

### 1.1 A brief intro to the reduced basis method

For many engineering purposes, one is interested in the evaluation of certain physical averages, or *outputs of interest*, defined as functionals of the solution to a partial differential equation (PDE) that describes the underlying physical problem. Often, the PDE is *parametrised* in the sense that the physical system is configured by a parameter vector  $\boldsymbol{\mu}$  governing e.g. boundary conditions, material properties, geometrical factors or loads. Given an output functional,  $l^{\text{out}}$ , which we shall assume to be linear and bounded, we write the output of interest as  $s(\boldsymbol{\mu}) = l^{\text{out}}(u(\boldsymbol{\mu}))$ , where  $u$  is the solution of the PDE corresponding to the parameter vector  $\boldsymbol{\mu}$ .

Consider the following situations:

- the parameter vector is unknown until just before the output is required (real-time control),
- the output is wanted for *many* different parameter vectors (optimisation, parameter estimation).

In the former context, computing a numerical approximation to  $u$  by a standard (say spectral element) method and then evaluating  $l^{\text{out}}(u)$  may not be possible within the short space of time available. In the latter context, using a standard method may be far too expensive in terms of computational cost. In both situations then, there is a premium on an alternative, faster method.

Let  $\mathcal{D} \subset \mathbb{R}^n$ ,  $n \in \mathbb{N}$ , be some admissible parameter space. Then, for a set of parameter vector samples  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_N \in \mathcal{D}$ , assume that the corresponding PDE

solutions  $u(\boldsymbol{\mu}_1), \dots, u(\boldsymbol{\mu}_N)$  are already available. We shall think of these functions as “snapshots” taken of  $u$  at different positions in parameter space. The reduced basis (RB) method now exploits the fact that if  $u$  is smooth in the parameter, i.e., the manifold  $\mathcal{M} = \{u(\boldsymbol{\mu}) : \boldsymbol{\mu} \in \mathcal{D}\}$  is smooth (which can be shown to be the case under certain hypothesis [19]), it should be possible to construct a good approximation  $u_N(\boldsymbol{\mu}) \approx u(\boldsymbol{\mu})$  as a linear combination of only a few snapshots for *any*  $\boldsymbol{\mu} \in \mathcal{D}$ , and hence arrive at a system of algebraic equations with only a few degrees of freedom. To find the optimal (in the “energy-norm” sense) linear combination, a standard Galerkin projection is used.

Hence, if the expense of precomputing “truth approximations” to the snapshots – to which a standard finite or spectral element method is employed – can be justified, the reduced basis method may, given any  $\boldsymbol{\mu} \in \mathcal{D}$ , drastically speed up the evaluation of the corresponding output of interest,  $s(\boldsymbol{\mu})$ .

## 1.2 Scope and overview

Our superior scope is to study and explore reduced basis techniques for partial differential equations. We focus on the empirical interpolation method and its applications to the reduced basis method for *non-affine* PDEs. We also develop an efficient method for the evaluation of flux type output functionals, i.e., if  $u_N$  is the numerical approximation to the field variable in a domain  $\Omega$ , then the output functional is on the form

$$(1.1) \quad I^{\text{out}}(u_N) = \int_{\Gamma} \frac{\partial u_N}{\partial n} ds,$$

where  $\Gamma \subset \partial\Omega$  and  $\frac{\partial u_N}{\partial n}$  denotes the outward normal derivative of  $u_N$ .

We start in Chapter 2 with preliminaries that we will make use of throughout the report. In particular, we consider certain theoretical aspects of a “Neumann-Dirichlet equivalence” that will prove very useful when evaluating flux-type output functionals.

Chapters 3, 4 and 5 are methodology chapters presenting the spectral element (SE), reduced basis (RB) and empirical interpolation (EI) methods, respectively. In Chapters 3 and 5, numerical examples are blended in to illustrate characteristic properties of the methods.

In Chapter 6, we exemplify the previously presented methodology and present a wide range of numerical results. The Neumann-Dirichlet equivalence and the SE, RB and EI methods will all be building blocks in an elaborate example

related to electrostatics. We also consider an inverse (parameter estimation) problem, making the most of the rapid input-output *forward* evaluation provided by these blocks put together.

From the lack of relevant literature, it seems that flux-type output functionals are rarely considered in the reduced basis context. In Chapter 6, our output of interest is a matrix of flux integrals, evaluated by way of the Neumann-Dirichlet equivalence from Chapter 2. As it turns out, this equivalence also readily allows for rapid evaluation of the output functionals. In fact, by also invoking empirical interpolation in our reduced basis approximation, our computations may be decoupled in an “offline-online” procedure, where the online stage – in which we, given any  $\boldsymbol{\mu} \in \mathcal{D}$ , compute the “RB-EI” output of interest – is very fast, and in particular independent of the computational complexity of the SE “truth approximations”.

### 1.3 A few remarks on notation

To reduce the possibility of future confusion, let us spend a few lines here clarifying some habits of notation.

Usually, vectors are denoted by an underline, e.g.  $\underline{v} \in \mathbb{R}^n$ . When we refer to the  $i$ 'th element of  $\underline{v}$ , we write  $(\underline{v})_i$ . As convenience requires, we will at times depart from this convention. One example, which we have already encountered, is the boldface parameter vector  $\boldsymbol{\mu}$ , which merits special attention when we are working with parametrised PDEs. Another example is the reduced basis solution coefficients, written as  $u_{N,1}, \dots, u_{N,N}$ , which are the elements of the solution vector  $\underline{u}_N$  (i.e., we omit the underline and parenthesis for the coefficients). At any rate, there should be no large risk of confusion whenever we deviate from our main rule.

As for vectors, we refer to element  $i, j$  of a matrix  $A$  by writing  $(A)_{ij}$ . Note that we do not denote matrices with an underline.

To particularly denote the parametric dependence of a function  $u = u(x, y)$  upon a parameter vector, we write  $u(x, y; \boldsymbol{\mu})$ , or simply  $u(\boldsymbol{\mu})$  when there is no need to emphasise the spatial dependence of  $u$ . For example, a reduced basis approximation to  $u(\boldsymbol{\mu})$  is written as  $u_N = u_N(x, y; \boldsymbol{\mu}) = u_N(\boldsymbol{\mu})$ , which is then not to be confused with the corresponding vector of solution coefficients, denoted by  $\underline{u}_N(\boldsymbol{\mu})$ .

On occasion, we use a right arrow to indicate asymptotic behaviour of a vari-

## 1. Introduction

---

able, e.g.  $a \rightarrow b$  or  $a \rightarrow \infty$ . This notation is not to be confused with the expression  $\boldsymbol{\mu} \rightarrow s(\boldsymbol{\mu})$ , denoting the evaluation of  $s$  for the parameter vector  $\boldsymbol{\mu}$ , nor with the left arrow used to denote assignment of values to variables in algorithm listings.

# Chapter 2

## Preliminaries

### 2.1 Weak form of the Poisson problem

We shall consider several Poisson problems in this report, written strongly as

$$(2.1) \quad -\kappa \Delta u(x, y) = f(x, y) \quad \text{in } \Omega \subset \mathbb{R}^2,$$

along with Dirichlet,

$$(2.2) \quad u(x, y) = g(x, y) \quad \text{on } \Gamma_D \subseteq \partial\Omega,$$

and Neumann,

$$(2.3) \quad \kappa \frac{\partial u}{\partial n}(x, y) = \rho(x, y) \quad \text{on } \Gamma_N \subseteq \partial\Omega,$$

boundary conditions, where  $\kappa \in \mathbb{R}$ ,  $\Gamma_D \neq \emptyset$  and  $\Gamma_D \cap \Gamma_N = \emptyset$ . We shall further assume that  $f \in L^2(\Omega)$ ,  $g \in L^2(\Gamma_D)$  and  $\rho \in L^2(\Gamma_N)$ , where we by  $L^2(\Pi)$  denote the usual space of square integrable functions over  $\Pi$ .

To derive the weak formulation of (2.1)–(2.3), we first define the spaces

$$(2.4) \quad X(\Omega) \stackrel{\text{def}}{=} \{v \in H^1(\Omega) : v|_{\Gamma_D} = 0\},$$

$$(2.5) \quad X^D(\Omega) \stackrel{\text{def}}{=} \{v \in H^1(\Omega) : v|_{\Gamma_D} = g\},$$

where  $H^1(\Omega)$  denotes the usual Sobolev space [3, 14, 25] of all functions belonging to  $L^2(\Omega)$  whose all first order derivatives also belong to  $L^2(\Omega)$ . Note that in the case of homogeneous Dirichlet boundary conditions ( $g \equiv 0$ ), we have

## 2. Preliminaries

---

$X = X^D$ . In the following, we may on occasion suppress the dependence of function spaces upon  $\Omega$  whenever no ambiguity may arise.

We arrive at the weak formulation of (2.1)–(2.3) by multiplying (2.1) with a test function  $v \in X$ , integrating and applying the Green's identity [25]

$$(2.6) \quad - \int_{\Omega} v \Delta u \, d\Omega = \int_{\Omega} (\nabla u)^T \nabla v \, d\Omega - \int_{\partial\Omega} v \frac{\partial u}{\partial n} \, ds.$$

Here,  $\frac{\partial u}{\partial n}$  denotes the outward normal derivative of  $u$  and  $ds$  the boundary measure on  $\partial\Omega$ . The weak formulation then reads: Find  $u \in X^D$  such that

$$(2.7) \quad \underbrace{\kappa \int_{\Omega} (\nabla u)^T \nabla v \, d\Omega}_{=a(u,v)} = \underbrace{\int_{\Omega} f v \, d\Omega + \int_{\Gamma_N} \rho v \, ds}_{=l(v)}, \quad \forall v \in X.$$

More often, we will write the problem abstractly in terms of the bilinear form  $a(\cdot, \cdot) : X^D \times X \rightarrow \mathbb{R}$  (linear in each of its arguments) and the linear functional  $l(\cdot) : X \rightarrow \mathbb{R}$ .

## 2.2 Norms and inner-products

The standard  $L^2$ ,  $H^1$  and  $X$  inner-products are, respectively,

$$(2.8) \quad (u, v)_{L^2} \stackrel{\text{def}}{=} \int_{\Omega} u^2 \, d\Omega, \quad (u, v)_{H^1} \stackrel{\text{def}}{=} \int_{\Omega} |\nabla u|^2 + u^2 \, d\Omega, \quad (u, v)_X \stackrel{\text{def}}{=} a(u, v),$$

with associated norms

$$(2.9) \quad \|u\|_{L^2} \stackrel{\text{def}}{=} \sqrt{(u, u)_{L^2}}, \quad \|u\|_{H^1} \stackrel{\text{def}}{=} \sqrt{(u, u)_{H^1}}, \quad \|u\|_{\mathcal{E}} \stackrel{\text{def}}{=} \sqrt{(u, u)_X},$$

where the latter is often referred to in literature as the “energy norm”. It can be shown [3] that whenever  $u$  is not a constant function, the energy-norm defined by  $a(\cdot, \cdot)$  in (2.7) and the  $H^1$ -norm are equivalent. On (rare) occasion, we will also make use of the infinity norm

$$(2.10) \quad \|v(x)\|_{L^\infty(\Omega)} \stackrel{\text{def}}{=} \sup_{x \in \Omega} |v(x)|.$$

Frequently, we shall invoke the Cauchy-Schwarz inequality [11]

$$(2.11) \quad |(v, w)| \leq ((v, v))^{1/2} ((w, w))^{1/2},$$

valid for all  $v, w \in X$  whenever  $(\cdot, \cdot)$  is an inner-product over  $X$ .



## 2.3 Gauss-Lobatto-Legendre Quadrature

Let  $\hat{\Omega} = (-1, 1) \times (-1, 1) \subset \mathbb{R}^2$ . For numerical evaluation of integrals, we use the Gauss-Lobatto-Legendre (GLL) quadrature formula

$$(2.12) \quad \int_{\hat{\Omega}} f(\xi, \eta) \, d\hat{\Omega} \approx \sum_{\alpha=0}^P \sum_{\beta=0}^P \rho_{\alpha} \rho_{\beta} f(\xi_{\alpha}, \xi_{\beta}),$$

where  $\xi_{\alpha}$ ,  $0 \leq \alpha \leq P$ , are the *GLL-nodes* and  $\rho_{\alpha}$ ,  $0 \leq \alpha \leq P$ , are the *GLL-weights* [25]. The formula (2.12) is exact for  $f(\xi, \eta) \in \mathbb{P}_{2P-1}(\hat{\Omega})$ , i.e., the space of polynomials of degree  $2P - 1$  or less in each spatial direction.

We may define Sobolev spaces  $H^{\sigma}$  for any real  $\sigma$  (see e.g. [3]). Suffice it here to say, somewhat loosely, that  $H^{\sigma}$  denotes the space of all square integrable functions whose all  $\sigma$ -order derivatives are also square integrable. Let  $\mathcal{I}_P u$  be the unique polynomial of degree  $P$  interpolating  $u$  in the tensorised GLL-nodes. If  $u \in H^{\sigma}$  with  $\sigma > 3/2$ , it can then be shown [3] that

$$(2.13) \quad \|v - \mathcal{I}_P v\|_{H^1} \leq c(\sigma) P^{1-\sigma},$$

where  $c(\sigma)$  is independent of  $P$ . Hence, for smooth integrands, GLL-quadrature will be “infinite-order” accurate.

## 2.4 Existence and uniqueness of a weak solution

First, let us formally establish the meaning of a few key terms. Let  $Z$  be a Hilbert space [11] equipped with the norm  $\|\cdot\|$ , and let  $a(\cdot, \cdot) : Z \times Z \rightarrow \mathbb{R}$  be a bilinear form. We then say that  $a$  is *coercive* provided there exists a “coercivity constant”

$$(2.14) \quad \alpha = \inf_{w \in Z} \frac{a(w, w)}{\|w\|^2} > 0,$$

and *continuous* if there exists a “continuity constant”

$$(2.15) \quad 0 < \gamma = \sup_{w \in Z} \sup_{v \in Z} \frac{a(w, v)}{\|w\| \|v\|} < \infty.$$

By  $Z'$ , we shall denote the dual space of  $Z$ , i.e. the space of bounded linear functionals over  $Z$ .

It is standard to ensure the existence of a solution to the problem (2.7) by appealing to the Lax-Milgram Theorem [3, 14, 25]. As we frequently rely on this theorem in the subsequent chapters, we include it (without proof) below.

**Theorem 2.1** (Lax-Milgram (existence and uniqueness)).

Let  $Z$  be a Hilbert space with norm  $\|\cdot\|$ ,  $a(\cdot, \cdot) : Z \times Z \rightarrow \mathbb{R}$  a coercive and continuous bilinear form and assume  $l \in Z'$ . Then, there exists a unique  $u \in Z$  such that

$$(2.16) \quad a(u, v) = l(v), \quad \forall v \in Z.$$

*Proof.* The reader is referred to e.g. [14] or [25] for a proof.  $\square$

For the spaces  $X$  and  $X^D$  defined in (2.4) and (2.5), respectively, note that Theorem 2.1 only explicitly covers homogeneous problems, i.e. problems in which the solution space,  $X^D$ , coincides with the space of test functions,  $X$ . This is automatically the case when either  $\Gamma_D = \emptyset$  or the boundary data  $g$  is identically zero (of course, for  $\Gamma_D = \emptyset$ , coercivity of  $a$  fails).

If we assume  $g$  not identically zero and let  $u^D$  denote the lifting of  $g$  into  $\Omega$  (that is to say,  $u^D$  is defined on all of  $\Omega$  and  $u^D|_{\partial\Omega} = g$ ), we may write  $u = u^D + u^0$  where  $u^0 \in X$ . Hence, we may instead consider the homogeneous problem: Find  $u^0 \in X$  such that

$$(2.17) \quad a(u^0, v) = -a(u^D, v) + l(v), \quad \forall v \in X,$$

which, according to Theorem 2.1, admits a unique solution whenever the right hand side,  $(-a(u^D, v) + l(v))$ , belongs to  $X'$ . It can be shown that this is the case if  $g$  is sufficiently regular, and in particular for  $g \in L^2(\Gamma_D)$ . In fact, the assumption  $g \in L^2(\Gamma_D)$  can be somewhat relaxed [25].

In the case that  $a(\cdot, \cdot)$  is symmetric, the Riesz Representation Theorem is sufficient to show the existence and uniqueness of a weak solution.

**Theorem 2.2** (Riesz Representation).

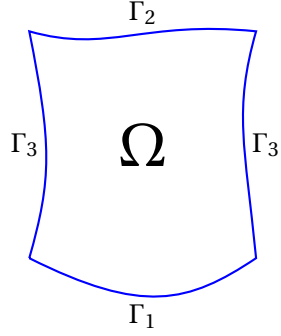
Let  $Z$  be a Hilbert space equipped with the inner-product  $(\cdot, \cdot)$  and let  $r \in Z'$ . Then there exists a unique element  $w \in Z$  such that

$$(2.18) \quad (w, v) = r(v), \quad \forall v \in Z.$$

*Proof.* The reader is referred to [14] for a proof.  $\square$

## 2.5 Evaluation of flux-type output functionals

In Chapter 6, we shall consider a flux-type output of interest from the reduced basis solution of a PDE. It seems from the lack of relevant literature that our par-



**Figure 2.1:**  $\Omega \subset \mathbb{R}^2$

particular type of output functional – allowed for by a “Neumann-Dirichlet equivalence”, discussed below – is previously unconsidered within the RB framework. In this section, we first describe the technique through an example and then state two lemmas making for some theoretical rigour.

Let  $\Omega$  be a two-dimensional domain with boundary  $\partial\Omega = \bar{\Gamma}_1 \cup \bar{\Gamma}_2 \cup \bar{\Gamma}_3$  (e.g. the domain depicted in Figure 2.1). For simplicity, we consider the Laplace problem

$$(2.19) \quad \begin{aligned} \Delta u &= 0 && \text{in } \Omega, \\ u &= g_1 && \text{on } \Gamma_1, \\ u &= g_2 && \text{on } \Gamma_2, \\ \frac{\partial u}{\partial n} &= 0 && \text{on } \Gamma_3, \end{aligned}$$

where  $g_1 \in L^2(\Gamma_1)$  and  $g_2 \in L^2(\Gamma_2)$ .

Now, let  $Y \subseteq H^1$  be a Hilbert space and define the function spaces

$$(2.20) \quad Z = \{v \in Y : v|_{\Gamma_1} = v|_{\Gamma_2} = 0\},$$

$$(2.21) \quad Z^D = \{v \in Y : v|_{\Gamma_1} = g_1, v|_{\Gamma_2} = g_2\},$$

and the bilinear form

$$(2.22) \quad a(u, v) = \int_{\Omega} \nabla v \cdot \nabla u \, d\Omega.$$

We then state the weak problem: Find  $u_w \in Z^D$  such that

$$(2.23) \quad a(u_w, v) = 0, \quad \forall v \in Z.$$

## 2. Preliminaries

---

An equivalent homogeneous formulation is: Find  $u_w^0 = u_w - u^D \in Z$  such that

$$(2.24) \quad a(u_w, v) \stackrel{\text{def}}{=} 0, \quad \forall v \in Z,$$

where  $u^D = u^{g_1} + u^{g_2} \in Z^D$  is a lifting of the Dirichlet data given by  $g_1$  and  $g_2$  to  $\Omega$  with  $u^{g_1}|_{\Gamma_2} = 0$  and  $u^{g_2}|_{\Gamma_1} = 0$ . We assume that the problem (2.23) admits a unique solution under the assumptions of Theorem 2.1.

The space  $Y$  in (2.20)–(2.21) is either equal to  $H^1$ , in which case (2.23) is the weak formulation of (2.19) and  $u_w$  is the corresponding weak solution (which we also refer to as the exact solution), or a discrete subspace of  $H^1$ , in which case  $u_w$  is a numerical approximation of  $u$ . When we intentionally refer to the exact solution (that is, when  $Y = H^1$ ), we omit the  $_w$  subscript.

Now, assume that our output of interest is the average outward flux through  $\Gamma_2$ , given by the flux integral

$$(2.25) \quad l^{\text{out}}(u_w) = - \int_{\Gamma_2} \frac{\partial u_w}{\partial n} ds,$$

where we without loss of generality have assumed that  $\Gamma_2$  is of unity length. The obvious thing to do next is to solve (2.23), evaluate  $l^{\text{out}}(u_w)$  by differentiating  $u_w$ , finding its normal derivative and then integrate along  $\Gamma_2$ . This procedure however, is subject to numerical differentiation and integration, and we may thus add to the output an additional numerical error. Furthermore, and particularly in the case of a curved boundary, the computation of  $\frac{\partial u_w}{\partial n}$  can be quite tedious to carry out as well.

We now look at an alternative to this “direct” output evaluation, and consider to this end the modified problem

$$(2.26) \quad \begin{aligned} \Delta u &= 0 && \text{in } \Omega, \\ u &= g_1 && \text{on } \Gamma_1, \\ \frac{\partial u}{\partial n} &= q && \text{on } \Gamma_2, \\ \frac{\partial u}{\partial n} &= 0 && \text{on } \Gamma_3, \end{aligned}$$

which is identical to the original problem (2.19) except for the replacement of the Dirichlet condition on  $\Gamma_2$  with a Neumann condition. Of course, if we choose  $q$  as the outward normal derivative across  $\Gamma_2$  of the solution to the original problem, the solution to (2.19) and (2.26) are identical.

Next, let

$$(2.27) \quad \tilde{Z} \stackrel{\text{def}}{=} Z \cup W,$$

where

$$(2.28) \quad W \stackrel{\text{def}}{=} \{v \in V : v|_{\Gamma_1} = 0, v|_{\Gamma_2} \neq 0\},$$

and  $V \subseteq H^1$  is a Hilbert space. Note that any  $v \in V$  which is identically zero on  $\Gamma_1$  and not identically zero on  $\Gamma_2$  is admissible. In the expanded space  $\tilde{Z}$ , we thus no longer enforce an essential Dirichlet condition on  $\Gamma_2$ . In general,  $V$  shall be different from  $Y$ .

We now state a (homogeneously written) weak problem: Find  $\tilde{u}_w^0 = \tilde{u}_w - u^{g_1} \in \tilde{Z}$  such that

$$(2.29) \quad a(\tilde{u}_w, v) = l(v), \quad \forall v \in \tilde{Z},$$

where

$$(2.30) \quad l(v) = \int_{\Gamma_2} \tilde{q} v \, ds.$$

Here,  $a(\cdot, \cdot)$  is given in (2.22) and  $u^{g_1}$  is a lifting of the boundary data  $g_1$  to  $\Omega$ . If  $\tilde{q} = q$  and  $Y = V = H^1$ , (2.29) is the weak formulation of (2.26) and  $\tilde{u}_w = u_w = u$  is the corresponding exact solution. On the other hand, if  $Y$  and  $V$  are discrete subspaces of  $H^1$ ,  $\tilde{u}_w$  is a discrete approximation to the solution of (2.26). In general,  $\tilde{q}$  shall be different from  $q$ .

Let us (until further notice) consider the particular case  $g_2 = 0 = u^{g_2}$ . We first state

**Lemma 2.1.**

*Assume that  $g_2$  is identically zero. If there exists a function  $\tilde{q} \in L^2(\Gamma_2)$  such that  $\tilde{u}_w|_{\Gamma_2} = u_w|_{\Gamma_2} (= 0)$ , then  $\tilde{u}_w = u_w$  in  $\Omega$ .*

*Proof.* By (2.29), it is clear that  $a(\tilde{u}_w, v) = 0$  for all  $v \in Z$  as  $Z \subset \tilde{Z}$  and the term on the right hand side vanishes for all  $v \in Z$ .

Next, as  $\tilde{u}_w|_{\Gamma_2} = \tilde{u}_w^0|_{\Gamma_2} = 0$ ,  $\tilde{u}_w^0$  can have no component in  $W$  by the definition (2.28) and hence  $\tilde{u}_w^0 \in Z$ . But then  $\tilde{u}_w = \tilde{u}_w^0 + u^{g_1} \in Z^D$ , and thus by uniqueness of the solution to the original problem (2.23),  $\tilde{u}_w = u_w$ .  $\square$

Consider the case in which  $Y = H^1$ , i.e., when  $u_w = u$  is the exact solution to (2.23). Under the assumptions of Lemma 2.1,  $u$  is also the exact solution to the modified problem (2.29) and hence  $\tilde{q} = \frac{\partial u}{\partial n} (= q)$  on  $\Gamma_2$  (to be precise,  $\frac{\partial u}{\partial n} = \tilde{q}$  almost everywhere on  $\Gamma_2$  [25]). In this case, we may evaluate the output of interest through the bilinear form as

$$(2.31) \quad l^{\text{out}}(u) = - \int_{\Gamma_2} \frac{\partial u}{\partial n} \, ds = - \int_{\Gamma_2} \tilde{q} \cdot 1 \, ds = -l(v^*) = -a(u, v^*),$$

## 2. Preliminaries

---

for any function  $v^\star \in V^\star \subset \tilde{Z}$ , where

$$(2.32) \quad V^\star \stackrel{\text{def}}{=} \{v \in W \subset \tilde{Z} : v|_{\Gamma_2} = 1\}.$$

Now assume  $Y$  to be a discrete space, and thus  $u_w$  is only a numerical approximation to  $u$ . Due to the weak imposition of the Neumann condition through the bilinear form in (2.29),  $\tilde{q}$  will in general be different from the outward normal derivative of  $u_w$  (or  $\tilde{u}_w$ ), i.e.  $\tilde{q} \neq \frac{\partial u_w}{\partial n}$  on  $\Gamma_2$ . However, as an approximation to the flux integral of (2.25), we may still choose to evaluate the numerical output of interest through the bilinear form as

$$(2.33) \quad l_{\text{app}}^{\text{out}}(u_w) = -a(u_w, v^\star), \quad v^\star \in V^\star.$$

Under the assumptions of Lemma 2.1, we have  $u_w = \tilde{u}_w$  and see that

$$(2.34) \quad \begin{aligned} l_{\text{app}}^{\text{out}}(u_w) &= -a(u_w, v^\star) = -a(\tilde{u}_w, v^\star) = - \int_{\Gamma_2} \tilde{q} v^\star \, ds \\ &\approx - \int_{\Gamma_2} \frac{\partial \tilde{u}_w}{\partial n} \, ds = - \int_{\Gamma_2} \frac{\partial u_w}{\partial n} \, ds, \end{aligned}$$

where we in the third equality invoke the problem definition (2.29). The immediately arising question is: if  $\tilde{q} \neq \frac{\partial u_w}{\partial n}$ , how well does  $l_{\text{app}}^{\text{out}}(u_w)$  approximate the output of interest  $l^{\text{out}}(u) = \int_{\Gamma_2} \frac{\partial u}{\partial n} \, ds$ ? Recalling that  $l^{\text{out}}(u) = a(u, v^\star)$ , it is straightforward to obtain the error estimate

$$(2.35) \quad \begin{aligned} |l_{\text{app}}^{\text{out}}(u_w) - l^{\text{out}}(u)| &= |a(u_w, v^\star) - a(u, v^\star)| \\ &= |a(u - u_w, v^\star)| \\ &\leq \|u - u_w\|_{\mathcal{E}} \|v^\star\|_{\mathcal{E}}, \end{aligned}$$

where we in the last step invoke the Cauchy-Schwarz inequality. Hence, if  $u_w \rightarrow u$  and  $\|v^\star\|_{\mathcal{E}}$  is bounded, the output error decays at least linearly with the energy norm of the error in the field variable.

We now make the obvious, but rather critical, note that we never actually solve the problem (2.29), and hence that the evaluation (2.33), and therein the choice of  $v^\star$ , does in fact define the “minimum required” expansion  $\tilde{X} \supset X$ , given in (2.27)–(2.28). In actual practice, we thus compute the discrete solution  $u_w$  of the original problem (2.23) and compute the numerical output from (2.33).

In general, the function  $\tilde{q}$ , and consequently  $l_{\text{app}}^{\text{out}}$ , will depend on the choice of  $v^\star$ . We derive the following result.

**Lemma 2.2.**

*For the problem (2.23) and under the assumptions of Lemma 2.1, two choices  $v_1^*$  and  $v_2^*$  for  $v^*$  in (2.33) are equivalent if  $w^* = v_1^* - v_2^* \in Z$ , and thus yield the same result when evaluating the output (2.33). On the contrary, this is not in general true if  $w^* \notin X$ .*

*Proof.* The proof is a simple consequence of the original problem definition (2.23), as

$$(2.36) \quad a(u_w, v_1^*) - a(u_w, v_2^*) = a(u_w, w^*) = 0 \quad \text{if } w^* \in Z.$$

On the other hand, if  $w^* \notin Z$ , then  $a(u_w, w^*)$  is not (in general) zero. Hence, the output evaluation (2.33) does in general depend on the choice of  $v^* \in V^*$ .  $\square$

Confined within the spectral element method framework, it seems natural to choose  $v^*$  as a polynomial of the same degree as the basis functions, and thus  $V = Y$ . Then we always have  $w^* \in Z$ , and thus every choice of  $v^*$  is equivalent by the previous lemma. In contrast, this is not true for the RB method. In Chapter 6 we shall consider examples of RB output evaluation with alternative choices of  $v^*$ .

When the boundary data on  $\Gamma_2$ ,  $g_2$ , is not identically zero, the arbitrariness of  $v^*$  “fails” at an earlier stage. Following the path of the proof of Lemma 2.1, we assume the existence of  $\tilde{q} \in L^2(\Gamma_2)$  such that  $\tilde{u}_w|_{\Gamma_2} = u_w|_{\Gamma_2} = g_2$ . Clearly, we still have  $a(\tilde{u}_w, v) = 0$  for all  $v \in Z$ . However,  $\tilde{u}_w$  must now have a nonzero component in  $W$  in order to represent the nonzero data on  $\Gamma_2$ . But then, we cannot in general conclude  $\tilde{u}_w \in Z^D$ , and thus  $\tilde{u}_w = u_w$ , unless  $V = Y$ .

Finally, we note that the error estimate (2.35) holds without regard to the choice of  $v^*$ .





# Chapter 3

## The Spectral Element (SE) method

The spectral element (SE) method, introduced in 1984 by A. T. Patera [21], is a member of the finite element methods family. By using only a few high-order elements, the method yields a high order of convergence while maintaining the geometrical flexibility of a standard low-order (e.g. linear or quadratic) finite element method.

In Chapter 6, we will employ the SE method in the construction of snapshots for a reduced basis approximation. However, as the method itself exhibits an interesting area of study, we will in this chapter explore it both theoretically and numerically.

We start by stating some of the basics of the spectral element method regarding formulation and numerical accuracy.

### 3.1 Spectral element discretisation

We consider the discretisation of a Poisson problem, abstractly written as: Find  $u \in X^D(\Omega)$  such that

$$(3.1) \quad a(u, v) = l(v), \quad \forall v \in X(\Omega).$$

Here,  $a(u, v) = \int_{\Omega} \nabla u \cdot \nabla v \, d\Omega$  is symmetric, continuous and coercive,  $l(v)$  is a bounded and linear functional (see Section 2.1) and  $X$  and  $X^D$  are defined in (2.4) and (2.5), respectively.

In the following, we assume that the physical domain  $\Omega$  admits partitioning

### 3. The Spectral Element (SE) method

---

into a connected union of  $K$  subdomains,

$$(3.2) \quad \Omega = \cup_{k=1}^K \overline{\Omega}_k,$$

where each of the subdomains  $\Omega_k$  are (possibly) deformed rectangles. As usual,  $\overline{\Omega} = \Omega \cup \partial\Omega$  is the closure of  $\Omega$ . We then define a *reference* or *computational domain*

$$(3.3) \quad \hat{\Omega} \stackrel{\text{def}}{=} (-1, 1) \times (-1, 1),$$

and assume that there from each physical subdomain exists a continuous one-to-one mapping onto the reference domain, formally

$$(3.4) \quad \mathcal{F}_k : \hat{\Omega} \rightarrow \Omega_k, \quad 1 \leq k \leq K.$$

The discretisation of (3.1) is carried out by choosing finite-dimensional substitutes  $X_{\mathcal{N}} \subset X$  for  $X$  and  $X_{\mathcal{N}}^D \subset X^D$  for  $X^D$ . We first define

$$(3.5) \quad \hat{v}_k(\xi, \eta) \stackrel{\text{def}}{=} v(x, y)|_{\Omega_k} \circ \mathcal{F}_k, \quad 1 \leq k \leq K.$$

The *reference function*  $\hat{v}_k$  is thus the restriction of  $v$  to  $\Omega_k$  mapped onto  $\hat{\Omega}$  through  $\mathcal{F}_k^{-1}$ . We shall refer to  $\xi, \eta$  as *reference variables*. As our discrete spaces, we then define

$$(3.6) \quad X_{\mathcal{N}}(\Omega) \stackrel{\text{def}}{=} \{v \in X(\Omega) : \hat{v}_k \in \mathbb{P}_P(\hat{\Omega}), 1 \leq k \leq K\},$$

$$(3.7) \quad X_{\mathcal{N}}^D(\Omega) \stackrel{\text{def}}{=} \{v \in X^D(\Omega) : \hat{v}_k \in \mathbb{P}_P(\hat{\Omega}), 1 \leq k \leq K\}.$$

The letter  $\mathcal{N}$  will serve two purposes. Firstly,  $\mathcal{N}$  denotes the SE discrete spaces and solutions and secondly,  $\mathcal{N}$  denotes the dimension of the discrete spaces, i.e

$$(3.8) \quad \mathcal{N} \stackrel{\text{def}}{=} \dim(X_{\mathcal{N}}^D) = \dim(X_{\mathcal{N}}),$$

and thus also the number of degrees of freedom associated with the numerical problem. In two dimensions,  $\mathcal{N} = \mathcal{O}(P^2)$ .

Finally, the discrete version of (3.1) can now be written abstractly as: Find  $u_{\mathcal{N}} \in X_{\mathcal{N}}^D$  such that

$$(3.9) \quad a(u_{\mathcal{N}}, v) = l(v), \quad \forall v \in X_{\mathcal{N}}.$$

## 3.2 Basis and algebraic formulation

The Lagrange interpolants through the  $P + 1$  GLL-nodes are given as

$$(3.10) \quad \ell_i(\xi) \stackrel{\text{def}}{=} \prod_{\substack{j=0 \\ j \neq i}}^P \frac{\xi - \xi_j}{\xi_i - \xi_j}, \quad i = 0, \dots, P.$$

As basis functions for the discrete spaces  $X_{\mathcal{N}}$  and  $X_{\mathcal{N}}^{\text{D}}$  we choose the tensorised Lagrange interpolants

$$(3.11) \quad \phi_{ij}(\xi, \eta) \stackrel{\text{def}}{=} \ell_i(\xi) \ell_j(\eta), \quad 0 \leq i, j \leq P.$$

In particular, we then have  $\phi_{ij}(\xi_m, \xi_n) = \delta_{im, jn}$  and (3.11) thus defines a nodal basis on the tensorised GLL-nodes. With Dirichlet boundary conditions, we simply omit the basis functions corresponding to the nodes on  $\Gamma_{\text{D}}$ . Note that the basis functions are defined on the reference domain  $\hat{\Omega}$ .

Through the mappings  $\mathcal{F}_k$ , we may write (3.9) in terms of the reference variables on  $\hat{\Omega}$ . Then writing the spectral element solution in terms of a linear combination of the basis functions, and letting the “ $\forall$ ” in (3.9) hold for each of the basis functions, we readily arrive at an algebraic formulation of the problem. So far, we have used a local (two-dimensional) numbering scheme, but if we instead assume a global (one-dimensional, e.g. lexographical) numbering of the unknowns and basis functions, we may write the system of algebraic equations as

$$(3.12) \quad A_{\mathcal{N}} \underline{u}_{\mathcal{N}} = \underline{l}_{\mathcal{N}},$$

where  $A_{\mathcal{N}}$  corresponds to the global spectral element stiffness matrix,  $\underline{l}_{\mathcal{N}}$  is the load vector and  $\underline{u}_{\mathcal{N}}$  is the vector of unknown coefficients.

## 3.3 Implementation notes and operation count

Let us briefly consider some key points regarding the implementation of our spectral element code. First, the integrals of the discrete formulation (3.9) are approximated numerically using GLL quadrature. We thus introduce a quadrature error to the numerical solution in addition to the already present approximation error. It is common to emphasise the quadrature evaluation of  $a$  and  $l$  by writing (say)  $a_{\mathcal{N}}$  and  $l_{\mathcal{N}}$ , respectively (as in e.g. [2, 3, 22]). Even though we

### 3. The Spectral Element (SE) method

---

always use GLL-quadrature to evaluate integrals, this notation is suppressed throughout this report.

Second, owing to (2.13) and the fact that our basis functions (3.11) are analytic ( $\sigma \rightarrow \infty$ ), we assume that integration over the  $(P+1)^2$  tensorised nodes yields a sufficiently accurate result when approximating the integrals in the stiffness matrix. Hence, only one set of GLL nodes needs to be defined. Together with the fact that the basis defined by (3.11) is nodal, this is particularly convenient when implementing a spectral element code [26].

Third, the global SE stiffness matrix,  $A_{\mathcal{N}}$ , is never explicitly formed, and a local, two-dimensional numbering scheme is used in actual practice. As a result, the effect of acting upon a vector with  $A_{\mathcal{N}}$  (operator evaluation) can be computed in only  $\mathcal{O}(P^3)$  floating point operations (flops) for  $\mathcal{O}(P^2)$  unknowns [26].

Finally, under the assumptions that  $a(\cdot, \cdot)$  is symmetric and coercive, the discrete operator is symmetric and positive definite. Hence, we can solve the system of algebraic equations by the conjugate gradient (CG) method [5], wherein operator evaluation and the Euclidean inner-product are the computational “kernels”. Assuming  $n_{\text{iter}}$  iterations of the CG algorithm, we obtain a spectral element solution in  $\mathcal{O}(n_{\text{iter}}P^3)$  flops. Considering the  $\mathcal{N} = \mathcal{O}(P^2)$  unknowns, this approach must be regarded as fairly efficient. In terms of memory requirements, we only need to store  $\mathcal{O}(P^2)$  floating points as the full stiffness matrix is never formed.

For comparison, note that a direct LU-solver, involving the explicit formation of the global stiffness matrix, would require a computational cost of  $\mathcal{O}(P^6)$  flops and a storage requirement of  $\mathcal{O}(P^4)$  floating points for  $\mathcal{O}(P^2)$  unknowns.

## 3.4 *A priori* error estimates

### 3.4.1 Error in the field variable

Define the *field error*

$$(3.13) \quad e_{\mathcal{N}} \stackrel{\text{def}}{=} u - u_{\mathcal{N}},$$

and assume for now that the SE solution  $u_{\mathcal{N}}$  is free of quadrature errors.

Since  $X_{\mathcal{N}} \subset X$ , we have the *Galerkin orthogonality* property

$$(3.14) \quad a(e_{\mathcal{N}}, v) = 0, \quad \forall v \in X_{\mathcal{N}}.$$

If we assume that  $a$  is continuous (with continuity constant  $\epsilon_2$ ), coercive (with coercivity constant  $\epsilon_1$ ) and symmetric, we get (with (3.14))

$$(3.15) \quad \epsilon_1 \|e_{\mathcal{N}}\|_{H^1}^2 \leq a(e_{\mathcal{N}}, e_{\mathcal{N}}) = a(e_{\mathcal{N}}, e_{\mathcal{N}}) + a(e_{\mathcal{N}}, u_{\mathcal{N}} - v) \\ = a(u - u_{\mathcal{N}}, u - v) \leq \epsilon_2 \|u - u_{\mathcal{N}}\|_{H^1} \|u - v\|_{H^1}, \quad \forall v \in X_{\mathcal{N}},$$

or more formally

**Theorem 3.1** (Céa's Lemma).

*Let  $u$  be the solution to (3.1) and  $u_{\mathcal{N}}$  the solution to (3.9). If  $a(\cdot, \cdot)$  is symmetric, coercive with coercivity constant  $\epsilon_1$  and continuous with continuity constant  $\epsilon_2$ , then*

$$(3.16) \quad \|u - u_{\mathcal{N}}\|_{H^1} \leq \frac{\epsilon_2}{\epsilon_1} \|u - v\|_{H^1}, \quad \forall v \in X_{\mathcal{N}}.$$

As a direct consequence of (3.14),  $u_{\mathcal{N}}$  will indeed be the best approximation to  $u$  when the error is measured in the energy norm and  $a$  is symmetric.

Let us for a while consider the case when  $\Omega = \hat{\Omega}$ . Let  $u \in H^\sigma(\hat{\Omega})$  with  $\sigma \geq 1$  and consider its best polynomial approximation  $u_P \in \mathbb{P}_P(\hat{\Omega})$  (when measured in the  $H^1$ -norm) defined by

$$(3.17) \quad u_P \stackrel{\text{def}}{=} \arg \min_{v \in \mathbb{P}_P(\hat{\Omega})} \|u - v\|_{H^1}^2.$$

It can be shown [4] that  $u_P$  satisfies

$$(3.18) \quad \|u - u_P\|_{H^1} \leq c(\sigma) P^{1-\sigma}$$

for every fixed  $\sigma \geq 1$ , where  $c(\sigma)$  is independent of  $P$ . If  $u$  is analytic, i.e. we may let  $\sigma \rightarrow \infty$ ,  $\|u - u_P\|_{H^1}$  decays faster than any algebraic power of  $1/P$  as  $P$  increases, and thus an exponential rate of convergence is achieved.

Now applying Theorem 3.1, there clearly exists  $\tilde{c} \in \mathbb{R}$  such that

$$(3.19) \quad \|u - u_{\mathcal{N}}\|_{H^1} \leq \tilde{c}(\sigma) P^{1-\sigma},$$

and the convergence rate is thus only limited by the regularity of  $u$ .

A similar result, which incorporates the contribution to the error from quadrature integration as well, holds for the Poisson problem in the multi-rectangle case. In particular, if the data  $f$  in (2.1) is smooth, as it will be in all problems considered in this report, the quadrature error vanishes exponentially fast and

### 3. The Spectral Element (SE) method

---

the convergence is again limited only by the regularity of  $u$  as in (3.19). Specifically, this result holds for  $\Omega$  defined by (3.2) as long as the boundary  $\partial\Omega$  is piecewise linear [2].

Two improved results, also incorporating quadrature errors, are also presented in [2] in the case of homogeneous Dirichlet boundary conditions. It is shown for smooth  $f$  and rectangular  $\Omega$  that

$$(3.20) \quad \|u - u_{\mathcal{N}}\|_{H^1} \leq c(\sigma)P^{-1-\sigma},$$

and if  $\Omega$  comprises one or more corners with interior angle equal to  $3\pi/2$  that

$$(3.21) \quad \|u - u_{\mathcal{N}}\|_{H^1} \leq c(\sigma)P^{1/3-\sigma},$$

where  $c(\sigma)$  is independent of  $P$ . Further, it is shown that the upper bound for the  $L^2$ -error is of order one better than that for the  $H^1$ -error. That is to say,

$$(3.22) \quad \|u - u_{\mathcal{N}}\|_{L^2} \leq c(\sigma)P^{-2-\sigma}$$

in the case of rectangular  $\Omega$ , and

$$(3.23) \quad \|u - u_{\mathcal{N}}\|_{L^2} \leq c(\sigma)P^{-2/3-\sigma}$$

in the case of non-convex corners.

The above results will be used for comparison when we analyse the error of numerical solutions.

#### 3.4.2 Error in the output of interest

Define the *output error*

$$(3.24) \quad e_{\mathcal{N}}^{\text{out}} \stackrel{\text{def}}{=} l^{\text{out}}(u) - l^{\text{out}}(u_{\mathcal{N}}),$$

where  $l^{\text{out}}$  is a linear and bounded output functional.

We now briefly describe the standard ‘‘Aubin-Nitsche’’ technique to argue that the output converges quadratically with the energy-norm error of the field variable. To this end, consider the *adjoint* problem of finding  $\psi \in X$  such that

$$(3.25) \quad a(v, \psi) = -l^{\text{out}}(v), \quad \forall v \in X.$$

The bilinear form  $a$  is the same as in the original, *primal* problem, and we have replaced the right hand side functional  $l$  with the output functional  $-l^{\text{out}}$  [20]. Correspondingly, we define the discrete problem: Find  $\psi_{\mathcal{N}} \in X_{\mathcal{N}}$  such that

$$(3.26) \quad a(v, \psi_{\mathcal{N}}) = -l^{\text{out}}(v), \quad \forall v \in X_{\mathcal{N}}.$$

Clearly, the field error  $e_{\mathcal{N}}$  defined in (3.13) belongs to  $X$ , and we thus have  $l^{\text{out}}(e_{\mathcal{N}}) = -a(e_{\mathcal{N}}, \psi)$  by (3.25). By also invoking Galerkin orthogonality and linearity of  $l^{\text{out}}$ , we get

$$(3.27) \quad e_{\mathcal{N}}^{\text{out}} = l^{\text{out}}(e_{\mathcal{N}}) = -a(e_{\mathcal{N}}, \psi) = -a(e_{\mathcal{N}}, \psi - \psi_{\mathcal{N}}).$$

Now, by the Cauchy-Schwarz inequality and continuity of  $a$  we arrive at

$$(3.28) \quad |e_{\mathcal{N}}^{\text{out}}| \leq a(e_{\mathcal{N}}, e_{\mathcal{N}})^{1/2} a(\psi - \psi_{\mathcal{N}}, \psi - \psi_{\mathcal{N}})^{1/2} \leq c \|e_{\mathcal{N}}\|_{H^1} \|\psi - \psi_{\mathcal{N}}\|_{H^1}$$

for some constant  $c > 0$ . If we again consider the case  $\Omega = \hat{\Omega}$ , assume  $u \in H^{\sigma_1}$  and  $\psi \in H^{\sigma_2}$ , then we may invoke (3.19) to get the estimate

$$(3.29) \quad |e_{\mathcal{N}}^{\text{out}}| \leq c_1(\sigma_1) c_2(\sigma_2) P^{1-\sigma_1} P^{1-\sigma_2}$$

for the error in the output. Assuming that  $\psi_{\mathcal{N}} \rightarrow \psi$  as fast as  $u_{\mathcal{N}} \rightarrow u$ , we have  $\sigma_1 = \sigma_2$  and the output converges with the same rate as the squared  $H^1$ -norm of the field variable.

In the particular case in which  $a$  is symmetric and  $l^{\text{out}} = l$ , the problem is said to be *compliant* [19]. This may for example be the case in a heat transfer (Poisson) problem where a heat flux is imposed on a part of the boundary and the output of interest is the average temperature over that particular boundary piece.

We note that with  $l^{\text{out}} = l$  in (3.26), and assuming that  $a$  is symmetric, the adjoint and primal problems coincide (apart from the minus sign), and the quadratic output convergence thus directly follows.

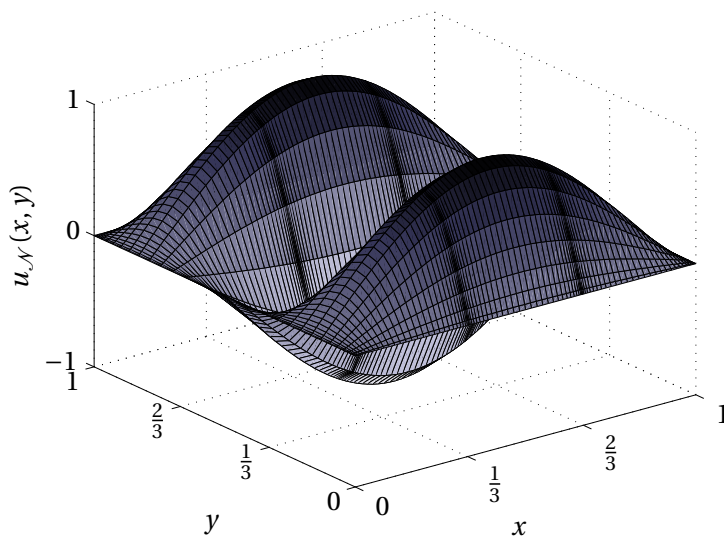
Finally, let us make an important notice. The estimate (3.29) holds generally, and depends merely on the definitions of  $a$ ,  $l$  and  $l^{\text{out}}$ . In particular, the result is independent of the discrete space  $X_{\mathcal{N}}$  and the choice of variational method. The recovery of the quadratic convergence (for the non-compliant case), is due to the richness of the discrete SE  $X_{\mathcal{N}}$ .

### 3.5 Numerical examples: Two model problems

To put the SE method to the test, we consider in this section two Poisson model problems to which the exact solutions are already explicitly known. Specifically, our examples fit within the assumptions of the *a priori* estimates from Section 3.4. Each example is accompanied by numerical results.

### 3. The Spectral Element (SE) method

---



**Figure 3.1:** A qualitative plot of the SE solution to the model problem (3.30). Three spectral elements are used. The polynomial degree is  $P = 22$ .

#### 3.5.1 A problem with known analytic solution

Consider the Poisson problem

$$(3.30) \quad \begin{aligned} -\Delta u &= f & \text{in } \Omega = (0, 1) \times (0, 1), \\ u &= 0 & \text{on } \partial\Omega, \end{aligned}$$

with  $f = 13 \sin(3\pi x) \sin(2\pi y)$ . The exact solution  $u = \sin(3\pi x) \sin(2\pi y)$  is analytic, so we expect exponential convergence when a spectral element solver is applied to solve the problem numerically.

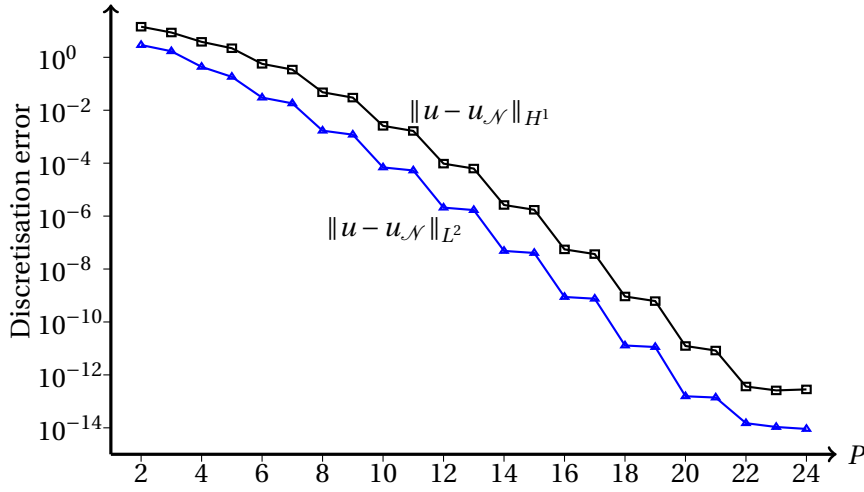
With, as is standard,  $H_0^1(\Omega)$  denoting the functions of  $H^1(\Omega)$  which vanish on  $\partial\Omega$ , the weak formulation of (3.30) reads: Find  $u \in X = H_0^1$  such that

$$(3.31) \quad \underbrace{\int_{\Omega} (\nabla u)^T \nabla v \, d\Omega}_{a(u, v)} = \underbrace{\int_{\Omega} f v \, d\Omega}_{l(v)}, \quad \forall v \in X.$$

It is readily shown that  $a$  is coercive and bounded, and that  $l$  is bounded. Thus, a unique solution to (3.31) exists by Theorem 2.1.

We now solve the problem numerically using the SE method. We use three spectral elements and a polynomial degree  $2 \leq P \leq 30$  on each element. A qualitative plot of the SE solution to (3.30) with a polynomial degree of  $P = 22$  is shown in Figure 3.1.





**Figure 3.2:** Errors as a function of the polynomial degree  $P$  for the SE solution to the model problem (3.30). Three spectral elements are used. The plot shows exponential convergence.

The convergence history of the numerical solution for increasing polynomial degree is shown in Figure 3.2. As expected from the *a priori* estimates of Section 3.4, both the energy-norm and  $L^2$ -norm error exhibit exponential convergence. The errors decay in steps, which seems reasonable due to the symmetry properties of the solution. When the polynomial degree  $P$  is increased by one, the decrease in the global error may depend strongly on whether  $P$  is odd or even. An odd polynomial may decrease the error quite a lot, whereas an even polynomial may be of little help – or vice versa.

### 3.5.2 A problem with known singularity solution

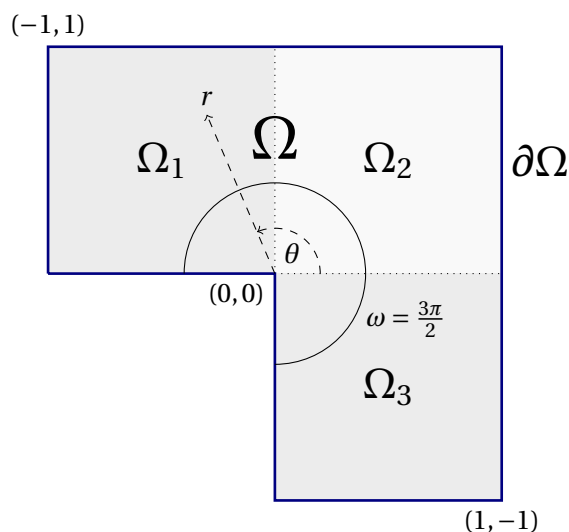
We shall now consider a problem for which the solution is known, but not analytic. First, we construct a singular solution to the Poisson problem. To this end, define

$$(3.32) \quad h(x) \stackrel{\text{def}}{=} \begin{cases} e^{-1/x}, & x > 0 \\ 0, & x \leq 0 \end{cases},$$

and then

$$(3.33) \quad c(r) \stackrel{\text{def}}{=} \frac{h\left(\frac{9}{10} - r\right)}{h\left(\frac{9}{10} - r\right) + h\left(r - \frac{1}{10}\right)}.$$

### 3. The Spectral Element (SE) method



**Figure 3.3:** The L-shaped domain  $\Omega$  for the model problem (3.35).

Now the “bump function”  $c(r) \in C^\infty(\mathbb{R})$  with the properties of being identically equal to 1 for  $r < 1/10$ , identically equal to 0 for  $r > 9/10$  and monotonically decreasing in between. For a proof of the analyticity of  $c(r)$ , confer [12, pages 49-51].

Consider the L-shaped domain depicted in Figure 3.3 and define in polar coordinates the function

$$(3.34) \quad u_s(r, \theta) \stackrel{\text{def}}{=} c(r)r^{\frac{2}{3}} \sin\left(\frac{2}{3}\left(\theta + \frac{\pi}{2}\right)\right).$$

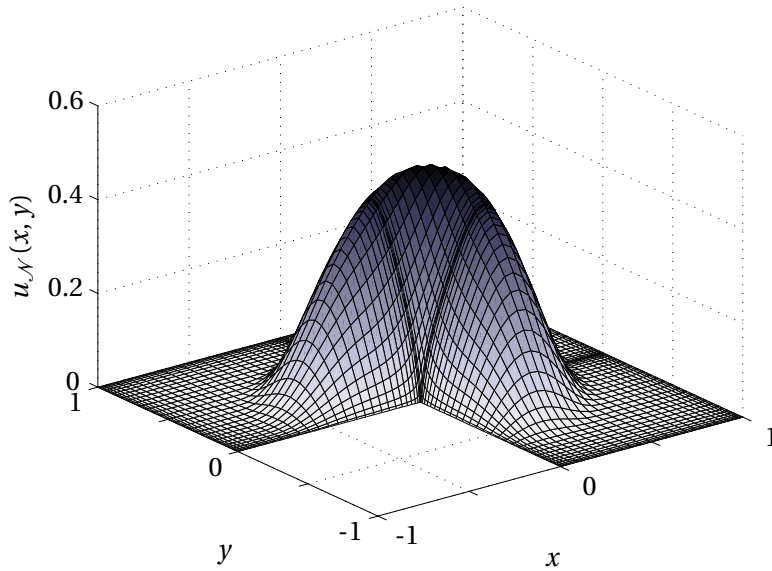
A singularity arises in  $\frac{\partial u_s}{\partial r}$  as  $r \rightarrow 0$ , so  $u_s$  is clearly not analytic. Plots of  $u_s$  are shown in Figures 3.4 and 3.5 (actually the plots show a spectral element solution of (3.35)). The singularity is readily seen on the contour plot (Figure 3.5), as the contour lines get increasingly close near the origin.

It can be shown [10, Lemma 2.4.1] that  $u_s$  satisfies

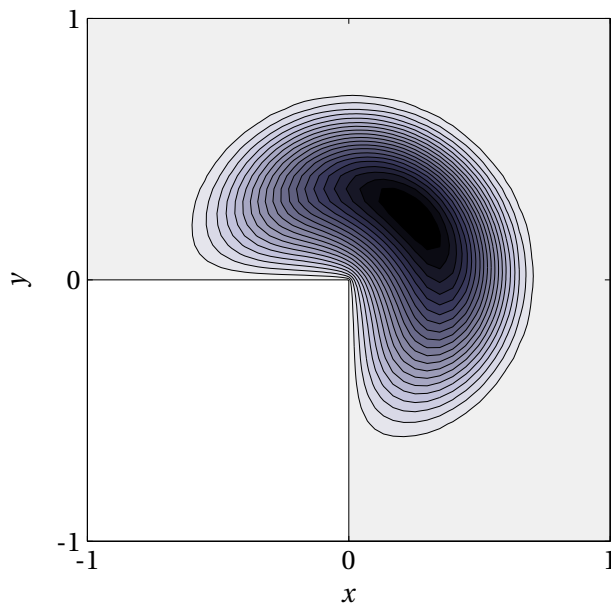
$$(3.35) \quad \begin{aligned} \Delta u_s &= F \quad \text{in } \Omega, \\ u_s &= 0 \quad \text{on } \partial\Omega \end{aligned}$$

with  $F$  analytic over  $\Omega$ . Given the angle  $\omega = 3\pi/2$  (Figure 3.3) and the expression given for  $u_s$  in (3.34), it can also be shown, without the explicit knowledge of  $u_s$  in (3.34), that  $u_s \in H^{5/3}(\Omega)$  [10, Remark 2.4.6].

The abstract and discrete formulations of (3.35) are readily derived. We then employ the SE method with three spectral elements as shown in Figure 3.3,



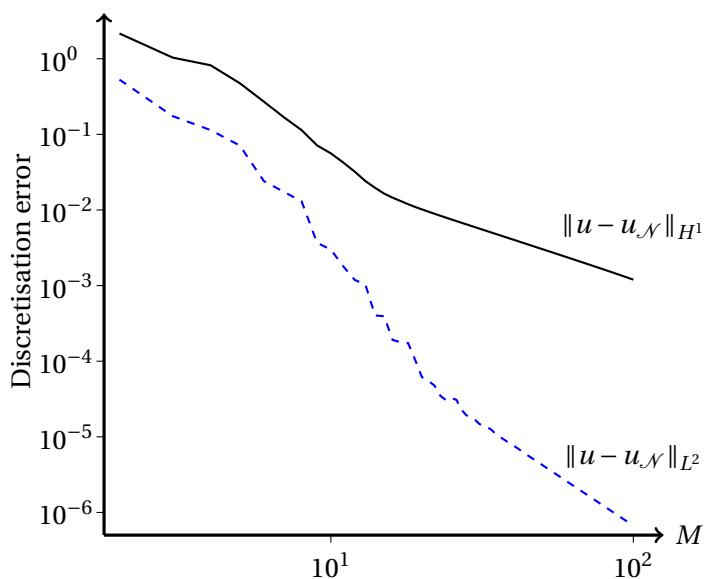
**Figure 3.4:** Surface plot of the SE solution to the model problem (3.35). The polynomial degree is  $M = 30$  and three spectral elements are used.



**Figure 3.5:** Contour plot of the SE solution to the model problem (3.35). The polynomial degree is  $P = 30$  and three spectral elements are used. The singularity is apparent in the plot through extremely dense contour lines close to the origin.

### 3. The Spectral Element (SE) method

---



**Figure 3.6:** Errors as a function of the polynomial degree  $2 \leq P \leq 100$  for the SE solution to (3.35). Three spectral elements are used. The plot shows algebraic convergence.

each with basis functions of degree  $P$  in each spatial direction. The numerical solution for  $P = 30$  is shown in Figures 3.4 (a “surface plot”) and 3.5 (a “contour plot”).

With  $u_s \in H^{5/3}(\Omega)$ , an algebraic convergence of order  $-2/3$  in the  $H^1$  norm is what we expect from the estimate (3.19) of Section 3.4. However, from the improved results (3.21) and (3.23) due to Bernardi and Maday [2], we expect a convergence rate of  $-4/3$  in the  $H^1$ -norm and  $-7/3$  in the  $L_2$ -norm, respectively.

Figure 3.6 shows the error in  $H^1$ -norm (solid) and  $L_2$ -norm (dashed) as functions of the polynomial degree  $2 \leq P \leq 100$ . For  $P \gtrsim 20$ , the decay seems to stabilise and convergence rates of  $-1.39$  and  $-2.66$  are achieved for the  $H^1$ - and  $L^2$ -norms, respectively. Clearly then, our findings are in agreement with the improved estimates (3.21) and (3.23) (as  $-1.3995 < -4/3$  and  $-2.6641 < -7/3$ ). We also conclude that for our problem, these error estimates are quite sharp.

# Chapter 4

## Reduced Basis (RB) Approximation

The reduced basis (RB) method may under certain assumptions provide a profound speedup to the computation of solutions of parametrised PDEs. In short, the idea is to tailor the approximation space specifically to the problem at hand by precomputing the solution for a small number of selected parameter vectors. These precomputed “snapshots” then span the discrete RB solution space, with the objective to reduce the number of required basis functions and thus the number of degrees of freedom.

In this chapter, we formulate the RB method for linear problems that are symmetric, coercive and parametrically affine. We describe a greedy algorithm for the selection of snapshots, as well as *a posteriori* error estimation and a very efficient offline-online computational procedure.

This chapter provides a short summary of the RB framework. For further reading, confer e.g. [19].

### 4.1 Formulation

#### 4.1.1 Parametric weak form

We first define our parameter space as

$$(4.1) \quad \mathcal{D} \subset \mathbb{R}^p, \quad p \in \mathbb{N}$$

for a modest number of parameters  $p$ . Each element of the parameter vector  $\boldsymbol{\mu} \in \mathcal{D}$  corresponds to a geometrical factor, a material property, a boundary

## 4. Reduced Basis (RB) Approximation

---

condition or something else that in some way configures the underlying partial differential equation.

Explicitly pointing out the parametrical dependence of the PDE, we now write  $a(\cdot, \cdot; \boldsymbol{\mu})$  and  $l(\cdot; \boldsymbol{\mu})$  for our bilinear form and linear functional, respectively. The parametric weak form becomes: For any  $\boldsymbol{\mu} \in \mathcal{D}$ , find  $u(\boldsymbol{\mu}) \in X$  such that

$$(4.2) \quad a(u, v; \boldsymbol{\mu}) = l(v; \boldsymbol{\mu}), \quad \forall v \in X,$$

and evaluate the output of interest

$$(4.3) \quad s(\boldsymbol{\mu}) = l^{\text{out}}(u(\boldsymbol{\mu})),$$

where  $l^{\text{out}}$  belongs to  $X'$ .

### 4.1.2 Norms and inner-products

The definitions of the now  $\boldsymbol{\mu}$ -dependent energy-norm and corresponding  $X$  inner-product follow naturally from the definitions in Section 2.2. Assuming that  $a(\cdot, \cdot; \boldsymbol{\mu})$  is coercive for all  $\boldsymbol{\mu} \in \mathcal{D}$  and for all  $v, w \in X$ , it is clear that

$$(4.4) \quad (v, w)_{\boldsymbol{\mu}} \stackrel{\text{def}}{=} a(v, w; \boldsymbol{\mu})$$

defines an inner-product on  $X$ . The associated parameter dependent energy norm is

$$(4.5) \quad \|w\|_{\mathcal{E}} \stackrel{\text{def}}{=} \sqrt{(w, w)_{\boldsymbol{\mu}}}.$$

Fixing  $\boldsymbol{\mu}_{\text{ref}} \in \mathcal{D}$  as a pre-defined *reference parameter vector*, we also define the parameter independent inner-product

$$(4.6) \quad (v, w)_{\boldsymbol{\mu}_{\text{ref}}} \stackrel{\text{def}}{=} a(v, w; \boldsymbol{\mu}_{\text{ref}}),$$

and the associated parameter independent norm

$$(4.7) \quad \|w\|_{\boldsymbol{\mu}_{\text{ref}}} \stackrel{\text{def}}{=} \sqrt{a(w, w; \boldsymbol{\mu}_{\text{ref}})}.$$

For all  $\boldsymbol{\mu} \in \mathcal{D}$ , we finally define the coercivity constant

$$(4.8) \quad 0 < \alpha(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \inf_{v \in X} \frac{a(v, v; \boldsymbol{\mu})}{a(v, v; \boldsymbol{\mu}_{\text{ref}})},$$

and the continuity constant

$$(4.9) \quad 0 < \gamma(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \sup_{v \in X} \sup_{w \in X} \frac{a(v, w; \boldsymbol{\mu})}{a(v, w; \boldsymbol{\mu}_{\text{ref}})} < \infty,$$

for  $a(\cdot, \cdot; \boldsymbol{\mu})$  with respect to the  $(\cdot, \cdot)_{\boldsymbol{\mu}_{\text{ref}}}$  inner-product.

### 4.1.3 Truth approximation

With  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_N \in \mathcal{D}$  we first define

$$(4.10) \quad S_N = \{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_N\}, \quad 1 \leq N \leq N_{\max}$$

as hierarchal sets of selected, distinct parameter vector samples. In general, the corresponding snapshots  $u(\boldsymbol{\mu}_1), \dots, u(\boldsymbol{\mu}_N)$  – which, in principle, shall span the RB approximation space – are unknown and in the need of numerical approximation. To this end, the spectral element method framework described in the previous chapter is employed (of course, standard finite element methods may equally well be considered).

For  $\boldsymbol{\mu} \in \mathcal{D}$ , let  $u_{\mathcal{N}_t}(\boldsymbol{\mu})$  be the spectral element *truth approximation* to  $u(\boldsymbol{\mu})$ , that is to say,  $u_{\mathcal{N}_t}(\boldsymbol{\mu}) \in X_{\mathcal{N}_t}$  subject to

$$(4.11) \quad a(u_{\mathcal{N}_t}, v; \boldsymbol{\mu}) = l(v; \boldsymbol{\mu}), \quad \forall v \in X_{\mathcal{N}_t}.$$

The number  $\mathcal{N}_t$  of degrees of freedom is assumed sufficiently large that the error  $\|u(\boldsymbol{\mu}) - u_{\mathcal{N}_t}(\boldsymbol{\mu})\|$  is “practically zero” for any desired norm. We shall refer to  $u_{\mathcal{N}_t}(\boldsymbol{\mu})$  as a truth approximation, as we assume that  $u_{\mathcal{N}_t}(\boldsymbol{\mu})$  and  $u(\boldsymbol{\mu})$  (the exact solution) are practically indistinguishable. Consequently,  $s_{\mathcal{N}_t}(\boldsymbol{\mu})$  and  $s(\boldsymbol{\mu})$  are also practically indistinguishable, and we shall refer to  $s_{\mathcal{N}_t}(\boldsymbol{\mu})$  as the truth output approximation.

Good truth approximations are important for two reasons. Firstly because they need to include as much information about the exact solution as possible to ensure rapid convergence of the RB solution, and secondly because we will estimate and measure (for error control and convergence analysis) the error in the RB solution relative to the truth approximations.

On occasion, we refer to  $u_{\mathcal{N}_t}$  and  $s_{\mathcal{N}_t}$  simply as the truth solution and output, respectively.

### 4.1.4 Discrete formulation

Given the parameter vector sets  $S_N$ ,  $1 \leq N \leq N_{\max}$ , we define the hierarchal reduced basis approximation spaces as

$$(4.12) \quad X_N = \text{span}\{u_{\mathcal{N}_t}(\boldsymbol{\mu}) : \boldsymbol{\mu} \in S_N\}, \quad 1 \leq N \leq N_{\max},$$

where we assume the  $u_{\mathcal{N}_t}(\boldsymbol{\mu})$ ,  $\boldsymbol{\mu} \in S_N$ , to be linearly independent. The discrete RB problem is then formulated as: For any given  $\boldsymbol{\mu} \in \mathcal{D}$ , find  $u_N(\boldsymbol{\mu}) \in X_N$  such

#### 4. Reduced Basis (RB) Approximation

---

that

$$(4.13) \quad a(u_N, v; \boldsymbol{\mu}) = l(v; \boldsymbol{\mu}), \quad \forall v \in X_N,$$

and evaluate the RB output

$$(4.14) \quad s_N(\boldsymbol{\mu}) = l^{\text{out}}(u_N(\boldsymbol{\mu})).$$

With the RB field error and RB output error defined as

$$(4.15) \quad e_N(\boldsymbol{\mu}) \stackrel{\text{def}}{=} u_{\mathcal{N}_t}(\boldsymbol{\mu}) - u_N(\boldsymbol{\mu})$$

and

$$(4.16) \quad e_N^{\text{out}}(\boldsymbol{\mu}) \stackrel{\text{def}}{=} s_{\mathcal{N}_t}(\boldsymbol{\mu}) - s_N(\boldsymbol{\mu}),$$

respectively, we state the following standard *a priori* optimality result.

**Theorem 4.1.**

*If  $a(\cdot, \cdot; \boldsymbol{\mu})$  is symmetric, coercive and continuous, then for  $\boldsymbol{\mu} \in \mathcal{D}$ ,  $u_{\mathcal{N}_t}(\boldsymbol{\mu})$  satisfying (4.11) and  $u_N(\boldsymbol{\mu})$  satisfying (4.13), we have*

$$(4.17) \quad \|u_{\mathcal{N}_t}(\boldsymbol{\mu}) - u_N(\boldsymbol{\mu})\|_{\mathcal{E}} = \inf_{w \in X_N} \|u_{\mathcal{N}_t}(\boldsymbol{\mu}) - w(\boldsymbol{\mu})\|_{\mathcal{E}}.$$

*Moreover, if  $l^{\text{out}} = l$  (i.e., the problem is compliant) we have*

$$(4.18) \quad e_N^{\text{out}}(\boldsymbol{\mu}) = \|e_N(\boldsymbol{\mu})\|_{\mathcal{E}}^2.$$

*Proof.* Equation (4.17) simply states that  $u_N$  is the optimal approximation in the energy-norm sense, which is a direct consequence of Galerkin orthogonality ( $X_N \subset X_{\mathcal{N}_t}$ ) and the symmetry of  $a$ .

For a compliant problem, we arrive at (4.18) by first using the fact that

$$(4.19) \quad e_N^{\text{out}}(\boldsymbol{\mu}) = l^{\text{out}}(e_N(\boldsymbol{\mu})) = l(e_N(\boldsymbol{\mu})) = a(u_{\mathcal{N}_t}, e_N; \boldsymbol{\mu}),$$

which follows from (4.11) as  $e_N \in X_{\mathcal{N}_t}$ . Then, by symmetry and Galerkin orthogonality, (4.18) follows.  $\square$

If we let  $\tilde{e}_N(\boldsymbol{\mu}) = u(\boldsymbol{\mu}) - u_N(\boldsymbol{\mu})$  and  $\tilde{e}_N^{\text{out}}(\boldsymbol{\mu}) = s(\boldsymbol{\mu}) - s_N(\boldsymbol{\mu})$ , we note that for any norm on  $X$  we have

$$(4.20) \quad \begin{aligned} \|\tilde{e}_N\| &= \|u - u_N\| = \|u - u_{\mathcal{N}_t} + u_{\mathcal{N}_t} - u_N\| \\ &\leq \|u - u_{\mathcal{N}_t}\| + \|u_{\mathcal{N}_t} - u_N\| = \|e_{\mathcal{N}_t}\| + \|e_N\| \end{aligned}$$



and

$$(4.21) \quad |\tilde{e}_N^{\text{out}}| = |s - s_N| = |s - s_{\mathcal{N}_t} + s_{\mathcal{N}_t} - s_N| \\ \leq |s - s_{\mathcal{N}_t}| + |s_{\mathcal{N}_t} - s_N| = |e_{\mathcal{N}_t}^{\text{out}}| + |e_N^{\text{out}}|$$

by the triangle inequality [11]. We thus readily see the importance of choosing  $X_{\mathcal{N}_t}$  rich enough that we for practical purposes can disregard the first term on the right-hand sides.

### 4.1.5 Algebraic formulation

So far in this report, we have used a local, two-dimensional numbering of our unknowns and basis functions. In the present context however, notation and analysis become easier and more general if we use a global, one-dimensional numbering scheme. In particular, we shall now write for our spectral element basis functions  $\psi_i(x, y)$  for  $1 \leq i \leq \mathcal{N}_t$ . Implementationwise, however, we still exploit the speedup and memory savings provided by a local, two-dimensional numbering scheme.

Every function  $v \in X_{\mathcal{N}_t}$  can now be written as

$$(4.22) \quad v = \sum_{i=1}^{\mathcal{N}_t} v_i \psi_i,$$

for some coefficients  $v_1, \dots, v_{\mathcal{N}_t}$ . Consequently, as  $a(\cdot, \cdot; \boldsymbol{\mu})$  is bilinear, we may for any  $v, w \in X_{\mathcal{N}_t}$  write

$$(4.23) \quad a(v, w; \boldsymbol{\mu}) = \sum_{i=1}^{\mathcal{N}_t} \sum_{j=1}^{\mathcal{N}_t} v_j w_i a(\psi_j, \psi_i; \boldsymbol{\mu}),$$

and

$$(4.24) \quad l(v; \boldsymbol{\mu}) = \sum_{i=1}^{\mathcal{N}_t} v_i l(\psi_i; \boldsymbol{\mu}).$$

To ensure numerical stability of our computations and to bound the condition number of the RB stiffness matrix [19], the snapshots  $u_{\mathcal{N}_t}(\boldsymbol{\mu}_1), \dots, u_{\mathcal{N}_t}(\boldsymbol{\mu}_N)$  are orthonormalised with respect to the  $\boldsymbol{\mu}_{\text{ref}}$ -inner-product to yield an orthonormal basis  $\zeta^1, \dots, \zeta^N$  for  $X_N$ . For  $1 \leq n \leq N$ , we thus have

$$(4.25) \quad \zeta^n(x, y) = \sum_{i=1}^{\mathcal{N}_t} \zeta_i^n \psi_i(x, y),$$

#### 4. Reduced Basis (RB) Approximation

---

for some coefficients  $\zeta_1^n, \dots, \zeta_{\mathcal{N}_t}^n$ . For  $1 \leq N \leq N_{\max}$ , we may now write the RB spaces as

$$(4.26) \quad X_N = \text{span}\{\zeta^1, \dots, \zeta^N\}.$$

The coefficients  $\zeta_1^n, \dots, \zeta_{\mathcal{N}_t}^n$ ,  $1 \leq n \leq N$ , are computed by way of a modified Gram-Schmidt orthogonalisation [5] of the original snapshot coefficients.

If we write the RB solution  $u_N(\boldsymbol{\mu})$  as

$$(4.27) \quad u_N(x, y; \boldsymbol{\mu}) = \sum_{n=1}^N u_{N,n}(\boldsymbol{\mu}) \zeta^n(x, y),$$

where  $u_{N,1}, \dots, u_{N,N}$  are the unknown coefficients, and let (4.13) be true for each of the basis functions of  $X_N$ , we arrive at the algebraic formulation

$$(4.28) \quad \sum_{n=1}^N u_{N,n}(\boldsymbol{\mu}) a(\zeta^n, \zeta^m; \boldsymbol{\mu}) = l(\zeta^m; \boldsymbol{\mu}), \quad 1 \leq m \leq N$$

of the reduced problem. The above system is equivalent of solving

$$(4.29) \quad A_N(\boldsymbol{\mu}) \underline{u}_N(\boldsymbol{\mu}) = \underline{l}_N(\boldsymbol{\mu}),$$

where  $\underline{u}_N(\boldsymbol{\mu}) = [u_{N,1}(\boldsymbol{\mu}), \dots, u_{N,N}(\boldsymbol{\mu})]^T \in \mathbb{R}^N$  is the unknown vector and the RB stiffness matrix  $A_N(\boldsymbol{\mu}) \in \mathbb{R}^{N \times N}$  and load vector  $\underline{l}_N(\boldsymbol{\mu}) \in \mathbb{R}^N$  are given element-wise as

$$(4.30) \quad (A_N(\boldsymbol{\mu}))_{mn} = a(\zeta^m, \zeta^n; \boldsymbol{\mu}) = \sum_{i=1}^{\mathcal{N}_t} \sum_{j=1}^{\mathcal{N}_t} \zeta_i^m \zeta_j^n a(\psi_i, \psi_j; \boldsymbol{\mu}), \quad 1 \leq m, n \leq N,$$

and

$$(4.31) \quad (\underline{l}_N(\boldsymbol{\mu}))_m = l(\zeta^m; \boldsymbol{\mu}) = \sum_{i=1}^{\mathcal{N}_t} \zeta_i^m l(\psi_i; \boldsymbol{\mu}), \quad 1 \leq m \leq N,$$

respectively, where  $a(\psi_i, \psi_j; \boldsymbol{\mu})$  and  $l(\psi_i; \boldsymbol{\mu})$  are the elements of the truth stiffness matrix,  $A_{\mathcal{N}_t}(\boldsymbol{\mu})$ , and load vector,  $l_{\mathcal{N}_t}(\boldsymbol{\mu})$ , respectively. If the basis functions for  $X_N$  are precomputed, we may then for any new  $\boldsymbol{\mu} \in \mathcal{D}$  assemble  $A_N(\boldsymbol{\mu})$  by (4.30) and  $\underline{l}_N(\boldsymbol{\mu})$  by (4.31), and then solve the presumably small system (4.29) to obtain a RB solution. Note that for  $1 \leq N \leq N_{\max}$ ,  $A_N$  is readily extracted from  $A_{N_{\max}}$  as the upper left  $N \times N$  submatrix, due to the fact that  $X_1 \subset \dots \subset X_{N_{\max}}$ .

We now make two remarks. Firstly, under the assumption that  $\mathcal{M} = \{u(\boldsymbol{\mu}) : \boldsymbol{\mu} \in \mathcal{D}\}$  is smooth, we expect that a good approximation to  $u(\boldsymbol{\mu})$  may be found in

$X_N$  for any new  $\boldsymbol{\mu} \in \mathcal{D}$ , even for fairly small  $N$ . Typically,  $N = \mathcal{O}(10)$ . In contrast, the number of unknowns associated with the spectral element system,  $\mathcal{N}_t$ , is likely to be very large as typical values for  $\mathcal{N}_t$  are  $\mathcal{O}(10^3)$  or even  $\mathcal{O}(10^4)$  (indeed,  $\mathcal{O}(10^5)$  or  $\mathcal{O}(10^6)$  degrees of freedom is not unusual in three-dimensional problems). Solving the system (4.29) then, is many times quicker than solving the system resulting from application of the spectral element method.

Secondly, note that the straightforward assembly of  $A_N(\boldsymbol{\mu})$  and  $l_N(\boldsymbol{\mu})$  in (4.30) and (4.31) require  $\mathcal{O}(N^2 \mathcal{N}_t^2)$  and  $\mathcal{O}(N \mathcal{N}_t)$  floating point operations, respectively. In actual practice however, we can often exploit either the sparsity of the truth stiffness matrix  $A_{\mathcal{N}_t}$  (the finite element case) or a local data representation (the spectral element case, see Section 3.3) to speed up the assembly of  $A_N$ . In any case, even when the basis functions are precomputed, the RB procedure still depends upon  $\mathcal{N}_t$  in terms of computational cost. Fortunately, for the important class of *affine* problems, it is possible to do all the  $\mathcal{N}_t$ -dependent computations in an offline-stage which is independent of the particular parameter vector for which to find the RB solution. We describe this computational strategy in the next section.

## 4.2 Offline-online procedure for affine problems

For affine problems, it is straightforward to develop an *offline-online* computational procedure which allows all the  $\mathcal{N}_t$ -complexity computations to be taken care of as part of the preprocessing stage. While the offline stage is computationally very expensive, the resulting online stage – in which we, given any  $\boldsymbol{\mu} \in \mathcal{D}$  compute the RB solution – is extremely fast. In particular, the online procedure is independent of  $\mathcal{N}_t$  in terms of computational cost.

An affine problem meets certain requirements on the bilinear form  $a(\cdot, \cdot; \boldsymbol{\mu})$  and linear functional  $l(\cdot; \boldsymbol{\mu})$ . First, the bilinear form may be written affinely in functions of the parameter vector as

$$(4.32) \quad a(\cdot, \cdot; \boldsymbol{\mu}) = \sum_{q=1}^{Q_a} \Theta_a^q(\boldsymbol{\mu}) a^q(\cdot, \cdot),$$

for a finite number  $Q_a$ . Here, the  $a^q(\cdot, \cdot)$  are parameter independent bilinear forms and the  $\Theta_a^q(\boldsymbol{\mu})$  are parameter dependent functions. We assume that  $Q_a$  is not too large, and that the parameter dependent functions are inexpensive to evaluate.

#### 4. Reduced Basis (RB) Approximation

---

Analogously, the linear functional may be written as

$$(4.33) \quad l(\cdot; \boldsymbol{\mu}) = \sum_{q=1}^{Q_l} \Theta_l^q(\boldsymbol{\mu}) l^q(\cdot),$$

for a modest number  $Q_l$ , where the  $l^q(\cdot)$  are parameter independent forms and the  $\Theta_l^q(\boldsymbol{\mu})$  are parameter dependent functions.

From (4.30) and (4.32) it is now evident that

$$(4.34) \quad \begin{aligned} (A_N(\boldsymbol{\mu}))_{mn} &= \sum_{i=1}^{\mathcal{N}_t} \sum_{j=1}^{\mathcal{N}_t} \zeta_i^m \zeta_j^n a(\psi_j, \psi_i; \boldsymbol{\mu}) \\ &= \sum_{i=1}^{\mathcal{N}_t} \sum_{j=1}^{\mathcal{N}_t} \zeta_i^m \zeta_j^n \left( \sum_{q=1}^{Q_a} \Theta_a^q(\boldsymbol{\mu}) a^q(\psi_j, \psi_i) \right) \\ &= \sum_{q=1}^{Q_a} \Theta_a^q(\boldsymbol{\mu}) \underbrace{\sum_{i=1}^{\mathcal{N}_t} \sum_{j=1}^{\mathcal{N}_t} \zeta_i^m \zeta_j^n a^q(\psi_j, \psi_i)}_{=(A_N^q)_{mn}}, \quad 1 \leq m, n \leq N, \end{aligned}$$

or equivalently,

$$(4.35) \quad A_N(\boldsymbol{\mu}) = \sum_{q=1}^{Q_a} \Theta_a^q(\boldsymbol{\mu}) A_N^q$$

for parameter independent  $N \times N$  matrices  $A_N^q$ . The approach for the right hand side is analogous. From (4.31) and (4.33) we get

$$(4.36) \quad \begin{aligned} (\underline{l}_N(\boldsymbol{\mu}))_m &= \sum_{i=1}^{\mathcal{N}_t} \zeta_i^m l(\psi_i; \boldsymbol{\mu}) \\ &= \sum_{i=1}^{\mathcal{N}_t} \zeta_i^m \left( \sum_{q=1}^{Q_l} \Theta_l^q(\boldsymbol{\mu}) l^q(\psi_i) \right) \\ &= \sum_{q=1}^{Q_l} \Theta_l^q(\boldsymbol{\mu}) \underbrace{\sum_{i=1}^{\mathcal{N}_t} \zeta_i^m l^q(\psi_i)}_{=(\underline{l}_N^q)_m}, \quad 1 \leq m \leq N, \end{aligned}$$

which is equivalent to

$$(4.37) \quad \underline{l}_N(\boldsymbol{\mu}) = \sum_{q=1}^{Q_l} \Theta_l^q(\boldsymbol{\mu}) \underline{l}_N^q$$

for parameter independent vectors  $\underline{l}_N^q \in \mathbb{R}^N$ .

The offline-online procedure is thus: Offline, we compute  $A_N^q$  for  $1 \leq q \leq Q_a$  and vectors  $\underline{l}_N^q$  for  $1 \leq q \leq Q_l$ . Online, we simply compute  $A_N(\boldsymbol{\mu})$  and  $\underline{l}_N(\boldsymbol{\mu})$  by (4.35) and (4.37) in  $\mathcal{O}(Q_a N^2)$  and  $\mathcal{O}(Q_l N)$  flops, respectively, and solve the small system (4.29). Presumably,  $N$  is small enough that a direct solver is the fastest alternative, and hence (4.29) is solved in  $\mathcal{O}(N^3)$  flops.

### 4.3 Snapshot sampling, a greedy algorithm

We now turn to the selection of parameter vectors for which to compute snapshot solutions. As our sets

$$(4.38) \quad S_N = \{\boldsymbol{\mu}_n\}_{n=1}^N, \quad 1 \leq N \leq N_{\max},$$

of parameter vectors are hierarchal, we construct  $S_{N+1}$  from  $S_N$  by including one new parameter vector at the time. Of course, we may simply choose the next vector manually, or randomly, from  $\mathcal{D}$ . Even though such an approach may, in fact, yield quite good results, it gives us little rigorous control of the approximation properties of  $X_N$  and hence of how good our reduced basis approximation will be.

An automatic and commonly used sampling procedure (in e.g. [22]) is a greedy algorithm. We start by assuming that

$$(4.39) \quad \Delta_N^{\text{out}}(\boldsymbol{\mu}) \geq |s_N(\boldsymbol{\mu}) - s_{\mathcal{N}_t}(\boldsymbol{\mu})|$$

is an upper bound for the RB output error for any  $\boldsymbol{\mu} \in \mathcal{D}$  (*a posteriori* error estimation is considered in the next section).

In short, the greedy algorithm searches  $\mathcal{D}$  for the parameter vector admitting the maximum value of  $\Delta_N(\boldsymbol{\mu})$ , includes it in the parameter vector set and expands the approximation space accordingly. Of course,  $\mathcal{D}$  consists of infinitely many points, so we have to settle with a surrogate “training sample”  $\Xi_{\text{train}} \subset \mathcal{D}$  of finite size. We assume that  $\Xi_{\text{train}}$  is fine enough that the behaviour of  $u(\boldsymbol{\mu})$  over  $\Xi_{\text{train}}$  is a good approximation of the behaviour of  $u(\boldsymbol{\mu})$  over  $\mathcal{D}$ .

Starting from a (say) randomly chosen initial parameter vector  $\boldsymbol{\mu}_1$  and corresponding approximation space  $X_1$ , the greedy algorithm constructs  $S_N$  and  $X_N$  for  $2 \leq N \leq N_{\max}$  in a way that at least in a sense is optimal. The procedure is listed below as Algorithm 4.1.

## 4. Reduced Basis (RB) Approximation

---



---

### Algorithm 4.1 Greedy parameter selection

---

Choose  $\boldsymbol{\mu}_1$ , compute  $u_{\mathcal{N}_t}(\boldsymbol{\mu}_1)$   
 $S_1 \leftarrow \{\boldsymbol{\mu}_1\}$ ,  $X_N \leftarrow u_{\mathcal{N}_t}(\boldsymbol{\mu}_1)$   
**for**  $2 \leq N \leq N_{\max}$  **do**  
     $\boldsymbol{\mu}_N \leftarrow \arg \max_{\boldsymbol{\mu} \in \Xi_{\text{train}}} \Delta_N^{\text{out}}(\boldsymbol{\mu})$   
     $S_N \leftarrow S_{N-1} \cup \boldsymbol{\mu}_N$   
     $X_N \leftarrow X_{N-1} \cup \text{span}\{u_{\mathcal{N}_t}(\boldsymbol{\mu}_N)\}$   
**end for**

---

## 4.4 *A posteriori* error estimation

### 4.4.1 Energy norm error bound

First, we assume that we may explicitly find a lower bound  $\alpha_{\text{LB}}(\boldsymbol{\mu})$  for the coercivity constant of  $a(\cdot, \cdot; \boldsymbol{\mu})$  with respect to the  $\boldsymbol{\mu}_{\text{ref}}$ -norm, i.e.,

$$(4.40) \quad \alpha_{\text{LB}}(\boldsymbol{\mu}) \leq \alpha(\boldsymbol{\mu}) = \inf_{v \in X_{\mathcal{N}_t}} \frac{a(v, v; \boldsymbol{\mu})}{a(v, v; \boldsymbol{\mu}_{\text{ref}})}.$$

Given the variational formulation (4.13) and a reduced basis solution  $u_N(\boldsymbol{\mu})$ , we also define the residual

$$(4.41) \quad r(v; \boldsymbol{\mu}) = l(v; \boldsymbol{\mu}) - a(u_N, v; \boldsymbol{\mu}), \quad \forall v \in X_{\mathcal{N}_t}.$$

As  $r$  belongs to  $X'_{\mathcal{N}_t}$  and  $a(\cdot, \cdot; \boldsymbol{\mu}_{\text{ref}})$  defines an inner-product on  $X_{\mathcal{N}_t}$ , we know by the Riesz Representation Theorem (Theorem 2.2) that there exists a unique  $\hat{e}(\boldsymbol{\mu}) \in X_{\mathcal{N}_t}$  such that

$$(4.42) \quad a(\hat{e}, v; \boldsymbol{\mu}_{\text{ref}}) = r(v; \boldsymbol{\mu}), \quad \forall v \in X_{\mathcal{N}_t}.$$

We then state formally

#### Theorem 4.2.

For  $\boldsymbol{\mu} \in \mathcal{D}$ , the *a posteriori* error estimator

$$(4.43) \quad \Delta_N(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \frac{\|\hat{e}(\boldsymbol{\mu})\|_{\boldsymbol{\mu}_{\text{ref}}}}{\sqrt{\alpha_{\text{LB}}(\boldsymbol{\mu})}} \geq \|e_N(\boldsymbol{\mu})\|_{\mathcal{E}},$$

where  $e_N(\boldsymbol{\mu}) = u_{\mathcal{N}_t}(\boldsymbol{\mu}) - u_N(\boldsymbol{\mu})$  and  $\alpha_{\text{LB}}(\boldsymbol{\mu})$  is the coercivity lower bound defined in (4.40).

*Proof.* With  $v = e_N$  in (4.42), the definition of the residual in (4.41) and the fact that  $e_N \in X_{\mathcal{N}_t}$ , we get

$$\begin{aligned}
 (4.44) \quad a(\hat{e}, e_N; \boldsymbol{\mu}_{\text{ref}}) &= r(e_N; \boldsymbol{\mu}) \\
 &= l(e_N; \boldsymbol{\mu}) - a(u_N, e_N; \boldsymbol{\mu}) \\
 &= a(u_{\mathcal{N}_t}, e_N; \boldsymbol{\mu}) - a(u_N, e_N; \boldsymbol{\mu}) \\
 &= a(e_N, e_N; \boldsymbol{\mu}).
 \end{aligned}$$

Then, by the Cauchy-Schwarz inequality,

$$(4.45) \quad a(e_N, e_N; \boldsymbol{\mu}) \leq \|\hat{e}\|_{\boldsymbol{\mu}_{\text{ref}}} \|e_N\|_{\boldsymbol{\mu}_{\text{ref}}}.$$

Finally, by letting  $v = e_N$  in Equation (4.40), we see that

$$(4.46) \quad \|e_N\|_{\boldsymbol{\mu}_{\text{ref}}} \leq \frac{(a(e_N, e_N; \boldsymbol{\mu}))^{1/2}}{\alpha_{\text{LB}}(\boldsymbol{\mu})^{1/2}}$$

and thus

$$(4.47) \quad a(e_N, e_N; \boldsymbol{\mu}) \leq \|\hat{e}\|_{\boldsymbol{\mu}_{\text{ref}}} \frac{(a(e_N, e_N; \boldsymbol{\mu}))^{1/2}}{\alpha_{\text{LB}}(\boldsymbol{\mu})^{1/2}},$$

from which (4.43) readily follows.  $\square$

For affine problems, an offline-online strategy for the computation of  $\|\hat{e}\|_{\boldsymbol{\mu}_{\text{ref}}}$  may also be developed, relying on the residual expansion (4.41) and the affine expansions (4.32) and (4.33). Moreover, even though finding a coercivity lower bound  $\alpha_{\text{LB}}$  is problem specific and not in general straightforward, a bound is readily found for problems that are *parametrically coercive* (the  $a^q(\cdot, \cdot)$  and  $\Theta_a^q$  are positive). An offline-online strategy for the *a posteriori* error estimator allows for *i*) rapid selection of snapshots during the greedy sampling procedure, and *ii*) rapid verification of the RB solution (and output, as we shall see below) in the RB online stage. For the details, the reader is referred to [19].

#### 4.4.2 Output error bounds

For a compliant problem, i.e.,  $a(\cdot, \cdot; \boldsymbol{\mu})$  is symmetric and  $l(\cdot; \boldsymbol{\mu}) = l^{\text{out}}(\cdot; \boldsymbol{\mu})$ , we can immediately deduce the “quadratic” bound

$$\begin{aligned}
 (4.48) \quad |s_{\mathcal{N}_t}(\boldsymbol{\mu}) - s_N(\boldsymbol{\mu})| &= l^{\text{out}}(e_N) \\
 &= a(e_N, e_N; \boldsymbol{\mu}) \\
 &\leq \Delta_N^2(\boldsymbol{\mu})
 \end{aligned}$$

#### 4. Reduced Basis (RB) Approximation

---

for the error in the RB output. The bound (4.48) follows directly from the problem definition (4.13), the fact that  $e_N \in X_{\mathcal{N}_t}$ , Galerkin orthogonality and symmetry of  $a$ .

Generally, and of particular interest for non-compliant problems, the ‘‘Aubin-Nitsche trick’’ considered in Section 3.4.2 for the SE method may be applied [23] to obtain

$$\begin{aligned}
 |s_{\mathcal{N}_t}(\boldsymbol{\mu}) - s_N(\boldsymbol{\mu})| &= l^{\text{out}}(e_N) \\
 (4.49) \quad &\leq (a(e_N, e_N; \boldsymbol{\mu}))^{1/2} (a(\psi_{\mathcal{N}_t} - \psi_N, \psi_{\mathcal{N}_t} - \psi_N; \boldsymbol{\mu}))^{1/2} \\
 &\leq \Delta_N(\boldsymbol{\mu}) (a(\psi_{\mathcal{N}_t} - \psi_N, \psi_{\mathcal{N}_t} - \psi_N; \boldsymbol{\mu}))^{1/2}.
 \end{aligned}$$

Here,  $\psi_{\mathcal{N}_t}$  is the truth solution to the dual problem: Find  $\psi_{\mathcal{N}_t}(\boldsymbol{\mu}) \in X_{\mathcal{N}_t}$  such that

$$(4.50) \quad a(v, \psi_{\mathcal{N}_t}; \boldsymbol{\mu}) = -l^{\text{out}}(v), \quad \forall v \in X_{\mathcal{N}_t},$$

and  $\psi_N$  is the corresponding RB approximation: Find  $\psi_N(\boldsymbol{\mu}) \in X_N$  subject to

$$(4.51) \quad a(v, \psi_N; \boldsymbol{\mu}) = -l^{\text{out}}(v), \quad \forall v \in X_N.$$

We recall that for the spectral element method, the ‘‘quadratic’’ convergence of non-compliant outputs is recovered as the SE approximation space,  $X_{\mathcal{N}}$ , is rich enough that the primal and dual errors decay equally fast. As the RB space  $X_N$  is specifically tailored to the *primal* problem however, we cannot expect to find  $\psi_N \in X_N$  such that  $\psi_{\mathcal{N}_t} - \psi_N$  is small, and hence the RB output error is proportional to  $\|e_N\|_{\mathcal{E}}$ .

However, as shown in [23], a quadratic RB output convergence rate may be recovered by including in  $X_N$  not only snapshots of  $u_{\mathcal{N}_t}$ , but of  $\psi_{\mathcal{N}_t}$  as well. It is also possible to consider the adjoint problem separately [20, 23], with its own RB approximation space. The latter approach is in general considered computationally advantageous [23], as two small systems can be solved in the RB online stage instead of a single, large system.



# Chapter 5

## Empirical Interpolation (EI)

To fully exploit the speedup offered by the RB method, an efficient offline-online computational procedure is required. In particular, the computational complexity in the online stage should be independent of  $\mathcal{N}_t$ . In Section 4.2 in the previous chapter, we described such a decoupling of computations for an affine problem. We recall the assumption of expansions of the bilinear form  $a(u, v; \boldsymbol{\mu})$  and linear functional  $l(v; \boldsymbol{\mu})$  on the form

$$(5.1) \quad a(\cdot, \cdot; \boldsymbol{\mu}) = \sum_{m=1}^{Q_a} \Theta_a^m(\boldsymbol{\mu}) a^m(u, v), \quad \text{and} \quad l(v; \boldsymbol{\mu}) = \sum_{n=1}^{Q_l} \Theta_l^n l^n(v),$$

for finite numbers  $Q_a$  and  $Q_l$ .

The bad news is that although many real-life problems admit an affine expansion, far from all do. With the Empirical Interpolation (EI) method however, it is possible to recover the online computational  $\mathcal{N}_t$ -independency.

The EI method was introduced in [1] and elaborated on in [9]. In this Chapter, we start by examining the method as it is presented in [1, 9], first theoretically and then numerically.

### 5.1 Motivation

To understand the implications of empirical interpolation to the RB method, consider as an example the one-dimensional bilinear form given by

$$(5.2) \quad a(u, v; \boldsymbol{\mu}) = \int_{-1}^1 \tau(x; \boldsymbol{\mu}) \frac{du}{dx} \frac{dv}{dx} dx.$$

## 5. Empirical Interpolation (EI)

---

For a general non-affine parameter dependent function  $\tau$  it is not possible to expand  $a$  on the form (5.1). The idea of the empirical interpolation method is to instead *approximate*  $\tau$  by an expansion

$$(5.3) \quad \tau(x; \boldsymbol{\mu}) \approx \tau_M(x; \boldsymbol{\mu}) = \sum_{m=1}^M \tilde{\varphi}_m(\boldsymbol{\mu}) \tau(x; \boldsymbol{\mu}_m^{\text{EI}}).$$

Here, the  $\tilde{\varphi}_m(\boldsymbol{\mu})$ , depend only on  $\boldsymbol{\mu}$ , and the functions  $\tau(x; \boldsymbol{\mu}_1^{\text{EI}}), \dots, \tau(x; \boldsymbol{\mu}_M^{\text{EI}})$  thus provide a set of “snapshots” of the non-affine function  $\tau$  taken at certain points in  $\mathcal{D}$ .

If we substitute  $\tau$  with  $\tau_M$  in (5.2), the resulting bilinear form

$$(5.4) \quad a_M(u, v; \boldsymbol{\mu}) = \int_{-1}^1 \tau_M \frac{du}{dx} \frac{dv}{dx} dx$$

is affine, and presumably a good approximation of the original bilinear form. Thus, for the bilinear form  $a_M$ , an offline-online decomposition is readily developed following the procedure from Section 4.2. On the downside, an additional error is introduced to the numerical solution due to the approximation (5.3), and hence  $M$  may be required quite large. The assumption then, is that  $\tau$  is *approximately affine* in the sense that  $\tau_M$  is a sufficiently good approximation to  $\tau$  for modest  $M$ .

Due to conditioning matters [1],  $\tau_M$  is in practice expressed as

$$(5.5) \quad \tau_M(x; \boldsymbol{\mu}) = \sum_{m=1}^M \varphi_m(\boldsymbol{\mu}) q_m(x),$$

where  $\{q_m\}_{m=1}^M$  is an orthogonal basis for

$$(5.6) \quad W_M \stackrel{\text{def}}{=} \text{span}\{\tau(x; \boldsymbol{\mu}_m^{\text{EI}})\}_{m=1}^M.$$

The EI method addresses *i*) the selection of parameter vectors  $\boldsymbol{\mu}_1^{\text{EI}}, \dots, \boldsymbol{\mu}_M^{\text{EI}}$  from which to compute the snapshots, and hence the basis functions  $q_m$ , and *ii*) the computation of the coefficients  $\varphi_m(\boldsymbol{\mu})$  given any new parameter vector  $\boldsymbol{\mu}$ . The former is taken care of by a greedy selection process, while the latter are computed by requiring that  $\tau_M$  be an interpolant of  $\tau$  at certain interpolation nodes  $\{t_1, \dots, t_M\}$  which are, in fact, also chosen in a greedy manner. Moreover, the stages *i*) and *ii*) conform to the offline-online procedure of the RB method; the construction of  $q_m$  ( $\boldsymbol{\mu}$ -independent) is done offline, while the evaluation of the  $\varphi(\boldsymbol{\mu})$  ( $\boldsymbol{\mu}$ -dependent) is quickly taken care of online.

The EI method has been successfully embedded in the RB framework for non-affine PDEs for a number of model problems, including problems that are time-dependent and non-linear, in e.g. [1, 9, 15]. Recently, it has also been employed in RB treatment of the Navier-Stokes equations with non-affine parametrical dependence [24].

## 5.2 Interpolation algorithm

---

### Algorithm 5.1 Empirical Interpolation

---

```

 $\boldsymbol{\mu}_1^{\text{EI}} \leftarrow \arg \max_{\boldsymbol{\mu} \in \Xi_{\text{train}}} \|\tau(\cdot; \boldsymbol{\mu})\|_{L^2}$ 
 $\xi_1(x) \leftarrow \tau(x; \boldsymbol{\mu}_1^{\text{EI}})$ 
 $t_1 \leftarrow \operatorname{argsup}_{x \in \Omega} |\xi_1(x)|$ 
 $q_1(x) \leftarrow \xi_1(x) / \xi_1(t_1)$ 
 $B_1 \leftarrow q_1(t_1)$ 
for  $2 \leq M \leq M_{\text{max}}$  do
     $\boldsymbol{\mu}_M^{\text{EI}} \leftarrow \arg \max_{\boldsymbol{\mu} \in \Xi_{\text{train}}} \varepsilon_{M-1}^*(\boldsymbol{\mu})$ 
     $\xi_M(x) \leftarrow \tau(x; \boldsymbol{\mu}_M^{\text{EI}})$ 
    Solve  $\sum_{j=1}^{M-1} \sigma_{M-1,j} (B_{M-1})_{ij} = \xi_M(t_i), \quad 1 \leq i \leq M-1$ 
     $r_M(x) \leftarrow \xi_M(x) - \sum_{j=1}^{M-1} \sigma_{M-1,j} q_j(x)$ 
     $t_M \leftarrow \operatorname{argsup}_{x \in \Omega} |r_M(x)|$ 
     $q_M(x) \leftarrow r_M(x) / r_M(t_M)$ 
     $(B_M)_{ij} \leftarrow q_j(t_i), \quad 1 \leq i, j \leq M$ 
end for
    
```

---

Let us now attend to the details of the EI algorithm, which is listed as Algorithm 5.1. Although we formally remain in the one-dimensional situation, meaning  $\Omega \subset \mathbb{R}^1$  is the physical domain of the problem, the extension to higher dimensions is obvious.

First, we introduce a finite training sample  $\Xi_{\text{train}}$  over the parameter space  $\mathcal{D}$ . We also define a matrix  $B_M \in \mathbb{R}^{M \times M}$  comprising the values of the basis functions  $q_1, \dots, q_M$  at the interpolation nodes  $t_1, \dots, t_M$ , given by

$$(5.7) \quad (B_M)_{ij} \stackrel{\text{def}}{=} q_j(t_i), \quad 1 \leq i, j \leq M.$$

The initial stage proceeds as follows. As the first parameter, we choose

$$(5.8) \quad \boldsymbol{\mu}_1^{\text{EI}} = \arg \max_{\boldsymbol{\mu} \in \Xi_{\text{train}}} \|\tau(\cdot; \boldsymbol{\mu})\|_{L^2}.$$

## 5. Empirical Interpolation (EI)

---

We then set  $\xi_1(x) = \tau(x; \boldsymbol{\mu}_1^{\text{EI}})$  and choose as the first interpolation node,

$$(5.9) \quad t_1 = \operatorname{argsup}_{x \in \Omega} |\xi_1(x)|.$$

Finally, we normalise the first basis function as  $q_1(x) = \xi_1(x)/\xi_1(t_1)$  and compute the (one element) matrix  $B_1$ .

After the initial operations for  $M = 1$ , we proceed in a similar manner for  $2 \leq M \leq M_{\max}$ . First, define the *projection error*

$$(5.10) \quad \varepsilon_M^*(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \inf_{z \in W_M} \|\tau(\cdot; \boldsymbol{\mu}) - z\|_{L^2}.$$

At each stage  $2 \leq M \leq M_{\max}$ , we then compute  $\varepsilon_M^*(\boldsymbol{\mu})$  for every  $\boldsymbol{\mu} \in \Xi_{\text{train}}$ , and make the greedy choice

$$(5.11) \quad \boldsymbol{\mu}_M^{\text{EI}} = \arg \max_{\boldsymbol{\mu} \in \Xi_{\text{train}}} \varepsilon_M^*, \quad \xi_M(x) = \tau(x; \boldsymbol{\mu}_M^{\text{EI}}).$$

To determine the best choice of new interpolation node  $t_M$ , we first interpolate  $\xi_M$  in the old interpolation nodes by solving the linear system

$$(5.12) \quad \sum_{j=1}^{M-1} \sigma_{M-1,j} \underbrace{q_j(t_i)}_{(B_{M-1})_{ij}} = \xi_M(t_i), \quad 1 \leq i \leq M-1,$$

for the unknown  $\sigma_{M-1,1}, \dots, \sigma_{M-1,M-1}$ . Then, we compute the *residual*

$$(5.13) \quad r_M(x) = \xi_M(x) - \sum_{j=1}^{M-1} \sigma_{M-1,j} q_j(x), \quad 1 \leq j \leq M-1.$$

For the next interpolation node, the greedy choice

$$(5.14) \quad t_M = \operatorname{argsup}_{x \in \Omega} |r_M(x)|$$

is made. Finally, we select as the next basis function the normalised residual

$$(5.15) \quad q_M(x) = \frac{r_M(x)}{r_M(t_M)},$$

and compute the matrix  $B_M$  of nodal values. As the residual vanishes at every old interpolation node by (5.12), i.e.,  $r_M(t_i) = 0$  for  $1 \leq i \leq M-1$ , and due to the “orthonormalisation” (5.15),  $B_M$  becomes lower triangular with unity diagonal.

We note that the greedy choice (5.11) of the parameter vectors, and consequently of the basis functions, ensures (in the greedy sense) optimal approximation spaces  $W_M = \text{span}\{q_m\}_{m=1}^M$ , whereas we hope the greedy choice of interpolation nodes (5.14) will ensure a small interpolation error *given*  $W_M$  as well.

From the basis for  $W_M$ , we can for any new  $\boldsymbol{\mu} \in \mathcal{D}$  construct a parametrically affine function  $\tau_M(x; \boldsymbol{\mu})$  given by

$$(5.16) \quad \tau_M(x; \boldsymbol{\mu}) = \sum_{i=1}^M \varphi_i(\boldsymbol{\mu}) q_i(x),$$

where the coefficients  $\varphi_i$ , only dependent upon  $\boldsymbol{\mu}$ , are determined such that  $\tau_M(x; \boldsymbol{\mu})$  is the interpolant of  $\tau(x; \boldsymbol{\mu})$  over  $\{t_1, \dots, t_M\}$ , i.e. by solving the system

$$(5.17) \quad \sum_{j=1}^M \varphi_j(\boldsymbol{\mu}) \underbrace{q_j(t_i)}_{(B_M)_{ij}} = \tau(t_i; \boldsymbol{\mu}), \quad 1 \leq i \leq M.$$

The empirical interpolation method admits an offline-online decomposition “by construction”. Offline, we run Algorithm 5.1 and invert the (presumably small) matrix  $B_M$ . Online, we may then for any  $\boldsymbol{\mu} \in \mathcal{D}$  quickly compute the coefficients  $\varphi_1(\boldsymbol{\mu}), \dots, \varphi_M(\boldsymbol{\mu})$  as the solution to the system (5.17), and then assemble the interpolant from (5.16).

### 5.2.1 A remark on practical implementation

In actual practice, we approximate all the functions involved in the above process by their polynomial interpolants of degree  $P_{\text{EI}}$  in the  $P_{\text{EI}} + 1$  GLL nodes, where  $P_{\text{EI}}$  is large enough that the interpolants are practically indistinguishable from the functions they approximate.

We are concerned with the computation of the projection error

$$(5.18) \quad \varepsilon_M^*(\boldsymbol{\mu}) = \inf_{z \in W_{M-1}} \|\tau(\cdot; \boldsymbol{\mu}) - z\|_{L^2}, \quad \boldsymbol{\mu} \in \Xi_{\text{train}}.$$

In the evaluation of the  $L^2$  norm, we approximate the integral  $\int_{\Omega} (\tau - z)^2 d\Omega$  with GLL quadrature, which in one dimension reads

$$(5.19) \quad \begin{aligned} \|\tau(\cdot; \boldsymbol{\mu}) - z\|_{L^2} &= \int_{-1}^1 (\tau - z)^2 d\Omega \\ &\approx \sum_{\alpha=0}^{P_t} \rho_{\alpha} (\tau(\xi_{\alpha}; \boldsymbol{\mu}) - z(\xi_{\alpha}))^2, \end{aligned}$$

## 5. Empirical Interpolation (EI)

---

where the  $\rho_\alpha$  are the GLL weights and the  $\xi_\alpha$  are the GLL nodes [25]. Now, let  $\Sigma \in \mathbb{R}^{(P_{\text{EI}}+1) \times (P_{\text{EI}}+1)}$  be a diagonal matrix comprising the GLL-weights, and write  $z \in W_M$  in terms of the basis functions as  $z = Q_M \underline{z}$  where  $Q_M = [\underline{q}_1, \dots, \underline{q}_M] \in \mathbb{R}^{(P_{\text{EI}}+1) \times M}$ , the  $\underline{q}_i$  are vectors comprising the values of  $q_i$  at the GLL-nodes and  $\underline{z} \in \mathbb{R}^M$  is a coefficient-vector. The minimisation problem (5.18) becomes

$$(5.20) \quad \varepsilon_M^*(\boldsymbol{\mu}) = \min_{\underline{z} \in \mathbb{R}^M} (Q_M \underline{z} - \underline{\tau})^T \Sigma (Q_M \underline{z} - \underline{\tau}),$$

where  $\underline{\tau}$  is a vector comprising the values of  $\tau$  in the GLL-nodes. But now the minimiser,  $\underline{z}^*$ , is simply the solution of the “weighted” normal equations

$$(5.21) \quad Q_M^T \Sigma Q_M \underline{z}^* = Q_M^T \Sigma \underline{\tau},$$

which may be solved using a standard Cholesky factorisation technique [5] (we note that since  $Q_M^T Q_M$  – the matrix arising from the standard least squares normal equations – is positive definite, then so is  $Q_M^T \Sigma Q_M$  due to the positive definiteness of  $\Sigma$ ).

### 5.3 Error estimation

The *projection error*

$$(5.22) \quad \varepsilon_M^*(\boldsymbol{\mu}) = \inf_{z \in W_M} \|\tau(\cdot; \boldsymbol{\mu}) - z\|_{L^2}$$

is (in general) different from the *interpolation error*

$$(5.23) \quad \varepsilon_M(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \|\tau(\cdot; \boldsymbol{\mu}) - \tau_M(\cdot; \boldsymbol{\mu})\|_{L^2}.$$

In [1, 9], it is shown that when measured in the  $L^\infty$ -norm, the interpolation error is only a constant away from the the projection error, when the number of interpolation nodes,  $M$  is fixed. To be precise, the result is

$$(5.24) \quad \|\tau(\cdot; \boldsymbol{\mu}) - \tau_M(\cdot; \boldsymbol{\mu})\|_{L^\infty} \leq (\Lambda(M) + 1) \left( \inf_{z \in W_M} \|\tau(\cdot; \boldsymbol{\mu}) - z\|_{L^\infty} \right),$$

where

$$(5.25) \quad \Lambda(M) \leq (2^M - 1).$$

We refer the reader to [1] or [9] for the proofs, and note that by the equivalence of norms in finite-dimensional spaces, a similar result will hold also for the  $L^2$ -norm of the interpolation error.

Increasing exponentially with  $M$ , the “constant”  $\Lambda(M)$  is undoubtedly too large to be of any practical value. It is established by numerical experiments however, that the bound (5.25) is likely to be far too conservative [9]. This claim is supported by the results presented in Section 5.4 in the sense that the ratio  $\varepsilon_M^*(\boldsymbol{\mu})/\varepsilon_M(\boldsymbol{\mu})$  increases only modestly with  $M$ . In the original paper [1], the greedy selection of parameters was based on the  $L^\infty$ -norm. However, [9] reports on only small differences in the interpolation error  $\|\tau(\cdot; \boldsymbol{\mu}) - \tau_M(\cdot; \boldsymbol{\mu})\|_{L^\infty}$ , whether the  $L^\infty$ -norm error or  $L^2$ -norm error is used as the criterion in the greedy selection process.

An *a posteriori* estimator for the interpolation error may also be developed, and incorporated in the reduced basis estimators, as shown in [1, 8, 9].

## 5.4 Numerical examples

### 5.4.1 A one-dimensional example

We now investigate a one-dimensional, single-parameter example. Consider the parametrically dependent function

$$(5.26) \quad f(x; \mu) = (x - 1)^\mu,$$

where  $x \in \Omega \stackrel{\text{def}}{=} (-1, 1)$  and  $\mu \in \mathcal{D} \stackrel{\text{def}}{=} [1, 4]$ .

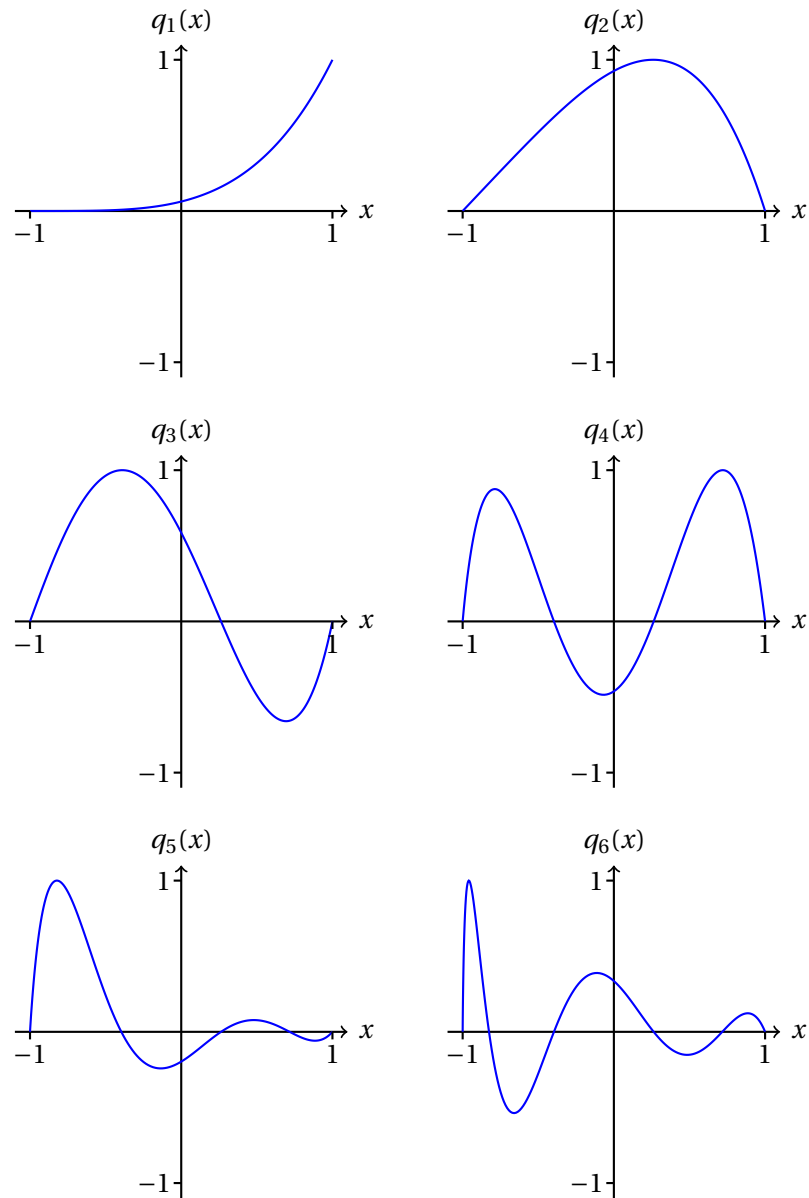
Clearly,  $f$  is non-affine in the parameter  $\mu$ . Moreover,  $f$  is weakly singular since  $f^{(n)}(-1) \rightarrow \infty$  when  $\mu$  is a non-integer and  $n > \mu$ . Here,  $f^{(n)}$  denotes the  $n$ 'th derivative of  $f$  with respect to  $x$ . Hence, by the polynomial interpolation error estimate (2.13), only algebraic convergence can be expected were we to interpolate  $f$  using standard GLL interpolation.

For the training sample  $\Xi_{\text{train}}$ , we choose a sample of 100 linearly distributed points over  $\mathcal{D}$ . We then apply the empirical interpolation method, as listed in Algorithm 5.1, to the function (5.26) with  $M_{\text{max}} = 15$  and  $P_{\text{EI}} = 200$ . We choose  $P_{\text{EI}}$  this large to make sure the singularity of  $f$  is well represented by its interpolant over the  $P_{\text{EI}} + 1$  GLL nodes (see Section 5.2.1).

A standard interpolation procedure makes use of basis functions that are common to every problem. In contrast, the empirical interpolation method tailors the approximation space  $W_M$  specifically to every new function (in the offline stage). As an example, we exhibit in Figure 5.1 the first six basis functions  $q_1, \dots, q_6$  for  $W_M$  when empirical interpolation is applied to the function (5.26).

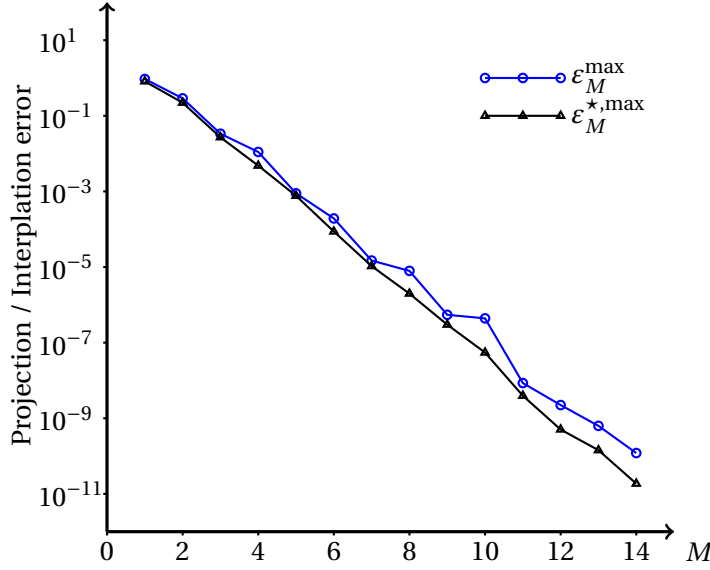
## 5. Empirical Interpolation (EI)

---



**Figure 5.1:** The first six basis functions for the one-dimensional empirical interpolation example.





**Figure 5.2:** Maximum projection error  $\varepsilon_M^{\star, \max}$  (triangles) and interpolation error  $\varepsilon_M^{\max}$  (circles) over  $\Xi_{\text{test}}$  as a function of  $M$  for the one-dimensional empirical interpolation example.

$M$	4	6	8	10	12	14
$\varepsilon_M^{r, \max}$	2.3628	2.2168	3.9920	8.3899	7.3410	7.3874

**Table 5.1:** Maximum ratio  $\varepsilon_M^{r, \max}$  of interpolation and projection errors over a test sample  $\Xi_{\text{test}}$ .

We note that the basis functions seem to be increasingly step in the vicinity of  $x = -1$ , suggesting that the EI method manages to capture the singular behaviour of  $f$ .

We next introduce a test sample  $\Xi_{\text{test}}$  of 100 randomly chosen points in  $\mathcal{D}$ . The plot in Figure 5.2 shows the maximum projection and interpolation errors over  $\Xi_{\text{test}}$ ,

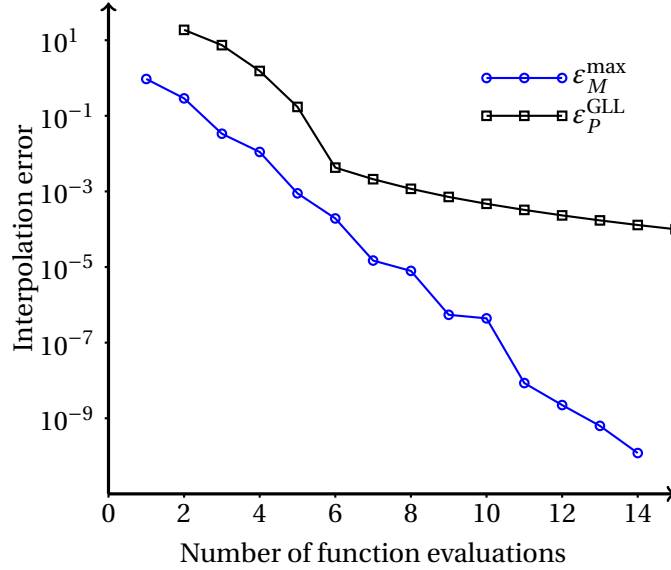
$$(5.27) \quad \varepsilon_M^{\star, \max} \stackrel{\text{def}}{=} \max_{\boldsymbol{\mu} \in \Xi_{\text{test}}} \varepsilon_M^{\star}(\boldsymbol{\mu}), \quad \varepsilon_M^{\max} \stackrel{\text{def}}{=} \max_{\boldsymbol{\mu} \in \Xi_{\text{test}}} \varepsilon_M(\boldsymbol{\mu}),$$

respectively, as  $M$  increases. We observe, in fact, an exponential rate of convergence with  $M$  for both  $\varepsilon_M^{\star, \max}$  and  $\varepsilon_M^{\max}$ . We also compute the maximum ratio

$$(5.28) \quad \varepsilon_M^{r, \max} \stackrel{\text{def}}{=} \max_{\boldsymbol{\mu} \in \Xi_{\text{test}}} \frac{\varepsilon_M(\boldsymbol{\mu})}{\varepsilon_M^{\star}(\boldsymbol{\mu})}$$

## 5. Empirical Interpolation (EI)

---



**Figure 5.3:** Interpolation error using GLL interpolation (squares) and empirical interpolation (circles) compared by the number of required function evaluations.

for  $M$  as listed in Table 5.1. As expected from Figure 5.2,  $\epsilon_M^{r,\max}$  increases only modestly with  $M$ .

Finally, let us consider briefly the alternative of approximating  $f$  by standard GLL interpolation. Letting  $\mathcal{I}_P f$  denote the polynomial interpolating  $f$  in the  $P + 1$  GLL nodes, we define

$$(5.29) \quad \epsilon_P^{\text{GLL}} = \max_{\boldsymbol{\mu} \in \Xi_{\text{test}}} \|f(\cdot; \boldsymbol{\mu}) - \mathcal{I}_P f(\cdot; \boldsymbol{\mu})\|_{L^2},$$

i.e., the maximum GLL interpolation error over the test sample, measured in the  $L^2$  norm. Figure 5.3 shows  $\epsilon_M$  and  $\epsilon_P^{\text{GLL}}$  as functions of the number  $M = P + 1$  of required evaluations of  $f$ . As expected from the estimate (2.13), GLL interpolation yields only algebraic asymptotic convergence due to the singularity at  $x = -1$ .

We conclude that for this particular problem, empirical interpolation is superior to standard polynomial interpolation. In fairness of GLL interpolation however, it should be noted that for a smooth  $f$  that are not too distorted we would expect more comparable results.

# Chapter 6

## A worked example: Electrostatics

In this chapter, we shall concretise the reduced basis and empirical interpolation frameworks described in the previous chapters by thoroughly examining a practical example. As we shall consider a two-dimensional electrostatics problem, the underlying partial differential equation is the Poisson equation.

Albeit a simple equation, we will encounter several difficulties along the way. For instance, the resulting bilinear form is non-affine due to non-affine terms in the parameter-dependent functions that describe the geometry of the physical domain. Moreover, we shall consider multiple non-compliant flux-type output functionals. Exploiting the Neumann-Dirichlet equivalence discussed in Section 2.5, we readily derive an efficient method for output evaluation. By also invoking the EI method, full  $\mathcal{N}_t$ -complexity decoupling of the offline and online RB stages is achieved.

We start with the derivation of a *forward* model, that is to say, the input-output relationship  $\boldsymbol{\mu} \rightarrow s(\boldsymbol{\mu})$ , where  $\boldsymbol{\mu}$  is a parameter vector and  $s$  is the output of interest. Herein lies the problem formulation, the construction of truth snapshot solutions for the underlying PDE, the reduced basis formulation and the associated *a posteriori* error estimation.

Lastly, we consider the corresponding *inverse* problem. Given a matrix  $s_{\text{obs}}$  of observed output data, we are seeking a parameter vector  $\boldsymbol{\mu}$  that are in compliance with the observations. However, as  $s_{\text{obs}}$  suffers from measurement noise, and the forward model from numerical (RB) error, it is in general impossible to find  $\boldsymbol{\mu}$  such that the corresponding output fits the observations exactly. In actual practice, we are instead trying to find  $\boldsymbol{\mu}$  such that  $s(\boldsymbol{\mu})$  is the minimiser of  $\|s_{\text{obs}} - s(\boldsymbol{\mu})\|$ , for some desired norm  $\|\cdot\|$ .

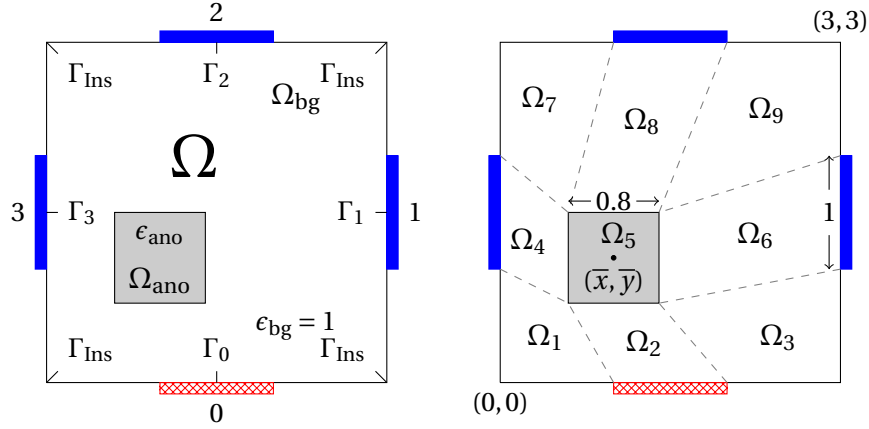


Figure 6.1: The domain  $\Omega$  with four electrodes.

## 6.1 Physical principles

### 6.1.1 The governing equation and boundary conditions

We consider the electrostatic potential,  $u$ , inside the two-dimensional domain  $\Omega = (0,3) \times (0,3)$  depicted in Figure 6.1. Attached to each of the edges of  $\Omega$  is an electrode on which a unity (red, crosshatched) or zero (blue) potential is imposed. The potential is equal to unity at one electrode only. Inside  $\Omega$  there is a background material  $\Omega_{\text{bg}}$  with electric permittivity  $\epsilon_{\text{bg}}$  and a small object, or anomaly,  $\Omega_{\text{ano}}$  with electric permittivity  $\epsilon_{\text{ano}} \neq \epsilon_{\text{bg}}$ . For simplicity, we assume that  $\epsilon_{\text{bg}} = 1$  and that  $\Omega$  and  $\Omega_{\text{ano}}$  are squares. Each of the edges of  $\Omega_{\text{ano}}$  is of length 0.8. The electrodes are of unity length, and are centred on each edge of  $\Omega$ .

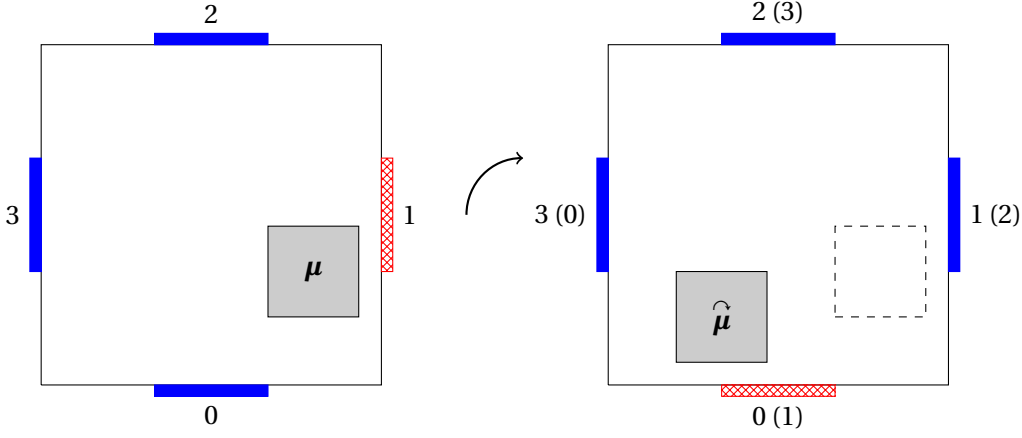
Denoting by  $\underline{E}$  the electrostatic field, we have by definition  $\underline{E} = -\nabla u$ . Inside  $\Omega$  we assume zero electric charge, and hence by Gauss' law  $\nabla \cdot \underline{E} = 0$ . For the potential  $u$  we thus arrive at the Laplace equation

$$(6.1) \quad -\Delta u = 0, \quad \text{in } \Omega_{\text{bg}} \cup \Omega_{\text{ano}}.$$

On the interior boundary  $\partial\Omega_{\text{ano}}$ , the “flux” continuity condition

$$(6.2) \quad \epsilon_{\text{bg}} \frac{\partial(u|_{\Omega_{\text{bg}}})}{\partial n} = -\epsilon_{\text{ano}} \frac{\partial(u|_{\Omega_{\text{ano}}})}{\partial n},$$

where  $\frac{\partial}{\partial n}$  denotes differentiation in the outward normal direction, holds. Finally, on the exterior boundary  $\Gamma_{\text{Ins}}$  between the electrodes, we assume electric



**Figure 6.2:** Rotational symmetry. By altering the configuration of the electrodes, the two systems are equivalent.

insulation, i.e.,

$$(6.3) \quad \frac{\partial u}{\partial n} = 0, \quad \text{on } \Gamma_{\text{Ins}}.$$

For a more thorough explanation of electrostatic boundary conditions, or indeed electromagnetic theory in general, the reader is referred to J. A. Stratton's classical textbook [27].

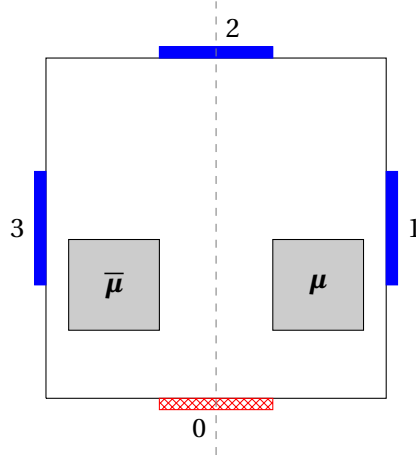
The capacitance corresponding to a red-blue pair of electrodes is given by  $C = Q/V$  where  $Q$  is the charge stored on either electrode and  $V = 1$  is the difference in potential between them. The stored charge on, say, electrode  $j$ , may be written as the flux integral  $Q = - \int_{\Gamma_j} \epsilon_{\text{bg}} \frac{\partial u}{\partial n} ds$ . We also introduce a parameter vector  $\boldsymbol{\mu} = (\bar{x}, \bar{y})$ , which configures the physical system in terms of the centre of  $\Omega_{\text{ano}}$  (see Figure 6.1). Now, we may measure and gather the capacitances in a symmetric *capacitance matrix*  $C = C(\boldsymbol{\mu}) \in \mathbb{R}^{4 \times 4}$  with elements

$$(6.4) \quad (C(\boldsymbol{\mu}))_{ij} = Q_j(\boldsymbol{\mu}) = - \int_{\Gamma_j} \frac{\partial u(\boldsymbol{\mu})}{\partial n} ds, \quad 0 \leq i, j \leq 3, i \neq j,$$

off the diagonal and zeroes on the diagonal.

### 6.1.2 Symmetry considerations

By switching which electrode on which to impose a unity potential, we may define four setups of the physical system that in principle may be handled sep-



**Figure 6.3:** Reflection of  $\mu$  across the vertical centreline.

arately from a mathematical point of view. By exploiting the rotational symmetry of the system however, it is sufficient to consider only the case with a unity potential on (say) electrode 0.

First, assume that  $\mu$  denotes the position of the anomaly as indicated to the left in Figure 6.2. Assuming that the electrostatic potential  $u(\mu)$  corresponding to a unity potential on electrode 0 is known, the elements  $C_{1j}(\mu)$ ,  $j = 1, 2, 3$ , of the capacitance matrix may be computed directly. As  $C$  is symmetric by definition, we only need to additionally consider the elements  $C_{12}(\mu)$ ,  $C_{13}(\mu)$  (unity potential on electrode 1) and  $C_{23}(\mu)$  (unity potential on electrode 2). Now, let  $\hat{\mu}$  denote a  $\pi/2$  clockwise rotation of  $\mu$  around the centre of  $\Omega$  as illustrated in Figure 6.2. It is then apparent that  $C_{12}(\mu) = C_{01}(\hat{\mu})$  and  $C_{13}(\mu) = C_{02}(\hat{\mu})$ . Similarly, we have  $C_{23}(\mu) = C_{03}(\hat{\mu})$ , where  $\hat{\mu}$  corresponds to a  $\pi/2$  counterclockwise rotation of  $\mu$  around the centre of  $\Omega$ . Surely, this symmetry argument relies on the fact that  $\Omega$  and  $\Omega_{\text{ano}}$  are square domains.

Our exploitation of symmetry doesn't stop here. Having reduced everything to only one red-blue configuration of the electrodes as described above, we now consider reflective symmetry across the vertical centreline, as depicted in Figure 6.3. Letting  $\bar{\mu}$  denote a reflection of  $\mu$  across the vertical centreline, we see that  $C_{03}(\mu) = C_{01}(\bar{\mu})$  and that  $C_{01}(\mu) = C_{03}(\bar{\mu})$ . Similarly,  $C_{02}(\mu) = C_{02}(\bar{\mu})$ . We will exploit these facts when constructing the discrete RB space, as we may in fact halve our parameter space  $\mathcal{D}$  and thus the number of required basis functions.

In addition, we note that the already mentioned symmetry of the capacitance

matrix  $C$  reflects the electrostatic reciprocity of the system – for a red-blue pair  $ij$  of electrodes,  $C_{ij}$  is determined by the *difference* in potential between the electrodes. We return (more mathematically) to this issue below.

## 6.2 RB treatment of the forward problem

### 6.2.1 Parametric weak form

As mentioned above, the parameter vector

$$(6.5) \quad \boldsymbol{\mu} \stackrel{\text{def}}{=} (\bar{x}, \bar{y}),$$

configures the physical system in terms of the spatial coordinates  $(\bar{x}, \bar{y})$  of the centre of the anomaly (see Figure 6.1). The parameter space  $\mathcal{D} \subset \mathbb{R}^2$  in which  $\boldsymbol{\mu}$  will reside is defined by

$$(6.6) \quad \mathcal{D} \stackrel{\text{def}}{=} (\bar{x}_{\min}, \bar{x}_{\max}) \times (\bar{y}_{\min}, \bar{y}_{\max}),$$

where

$$(6.7) \quad \begin{aligned} \bar{x}_{\min} &= 3/2, \\ \bar{y}_{\min} &= 1, \\ \bar{x}_{\max} &= \bar{y}_{\max} = 2. \end{aligned}$$

Note that we implicitly include the “mirrored” parameter space

$$(6.8) \quad \overline{\mathcal{D}} \stackrel{\text{def}}{=} [1, 3/2] \times [1, 2]$$

as well, due to the flip-symmetry described above.

From the rotational symmetry argument, we shall only consider the case with an imposed unity potential on electrode 0 (and a zero potential on electrodes 1, 2 and 3). Hence, we only need to define and solve problems corresponding to a single weak formulation.

Note that the parameter vector  $\boldsymbol{\mu} = (\bar{x}, \bar{y})$ , does in fact change the shape of  $\Omega_{\text{bg}}$  and the position of  $\Omega_{\text{ano}}$ . We may thus write  $\Omega_{\text{bg}} = \Omega_{\text{bg}}^{\boldsymbol{\mu}}$  and  $\Omega_{\text{ano}} = \Omega_{\text{ano}}^{\boldsymbol{\mu}}$  to emphasise the  $\boldsymbol{\mu}$ -dependency of the physical domains. As our bilinear form, we then define

$$(6.9) \quad a(u, v; \boldsymbol{\mu}) \stackrel{\text{def}}{=} \int_{\Omega_{\text{bg}}^{\boldsymbol{\mu}}} \epsilon_{\text{bg}} \nabla u \cdot \nabla v \, d\Omega_{\text{bg}}^{\boldsymbol{\mu}} + \int_{\Omega_{\text{ano}}^{\boldsymbol{\mu}}} \epsilon_{\text{ano}} \nabla u \cdot \nabla v \, d\Omega_{\text{ano}}^{\boldsymbol{\mu}},$$

## 6. A worked example: Electrostatics

---

and as our test and approximation spaces

$$(6.10) \quad X \stackrel{\text{def}}{=} \{v \in H^1(\Omega) : v|_{\{\Gamma_i\}_{i=0}^3} = 0\}$$

$$(6.11) \quad X^D \stackrel{\text{def}}{=} \{v \in H^1(\Omega) : v|_{\Gamma_0} = 1, v|_{\{\Gamma_i\}_{i=1}^3} = 0\},$$

respectively. Due to typesetting convenience, we suppress the  $\boldsymbol{\mu}$ -dependency of the physical domains in what follows. The parametric weak form of the Laplace equation (6.1) with the Neumann conditions (6.2) and (6.3) and Dirichlet data

$$(6.12) \quad u|_{\Gamma_i} = \begin{cases} 1, & \text{if } i = 0, \\ 0, & \text{if } i = 1, 2, 3, \end{cases}$$

becomes: Given  $\boldsymbol{\mu} \in \mathcal{D}$ , find  $u(\boldsymbol{\mu}) \in X^D$  such that

$$(6.13) \quad a(u, v; \boldsymbol{\mu}) = 0, \quad \forall v \in X.$$

The electric insulation on  $\Gamma_{\text{Ins}}$  and the interior flux continuity condition (6.2) are naturally taken care of by the weak form (6.13), whereas the Dirichlet conditions (6.12) and the global  $C^0$  continuity of  $u$  is ensured by the choice of  $X$  and  $X^D$ .

Finally, as our reference parameter vector, we choose the centre of  $\mathcal{D}$ , namely  $\boldsymbol{\mu}_{\text{ref}} \stackrel{\text{def}}{=} (1.75, 1.5)$ .

### 6.2.2 Spectral element truth approximation

#### Discrete formulation

For a spectral element approximation to  $u$ , we recall from Chapter 3 the assumption of partitioning of the physical domain into a finite number of deformed rectangles, which in turn may be mapped onto the reference domain  $\hat{\Omega}$ . To this end, we define  $\Omega = \cup_{i=1}^9 \bar{\Omega}_i$  as depicted (to the right) in Figure 6.1. We further define continuous one to one mappings

$$(6.14) \quad \mathcal{F}_k : \hat{\Omega} \rightarrow \Omega_k, \quad 1 \leq k \leq 9,$$

mapping the reference domain onto each of the physical subdomains. The mapping for  $\Omega_1$  is explicitly considered below.



Given the partitioning of  $\Omega$  and the associated mappings, the discrete spectral element “truth” spaces  $X_{\mathcal{N}_t}^D$  and  $X_{\mathcal{N}_t}$  are properly defined by

$$(6.15) \quad \begin{aligned} X_{\mathcal{N}_t}^D(\Omega) &\stackrel{\text{def}}{=} \{v \in V(\Omega) : v|_{\{\Gamma_i\}_{i=0}^3} = 0\} \\ X_{\mathcal{N}_t}(\Omega) &\stackrel{\text{def}}{=} \{v \in V(\Omega) : v|_{\Gamma_1} = 0, v|_{\{\Gamma_i\}_{i=1}^3} = 0\}, \end{aligned}$$

where

$$(6.16) \quad V(\Omega) \stackrel{\text{def}}{=} \{v \in H^1(\Omega) : v|_{\Omega_k} \circ \mathcal{F}_k \in \mathbb{P}_{P_t}(\hat{\Omega}), 1 \leq k \leq 9\}.$$

We shall assume that the truth spaces are sufficiently rich when we use basis functions of polynomial degree  $P_t \stackrel{\text{def}}{=} 30$ .

Finally, for the spectral element truth approximation  $u_{\mathcal{N}_t}(\boldsymbol{\mu})$  to the electric potential  $u(\boldsymbol{\mu})$ , we are thus seeking  $u_{\mathcal{N}_t}(\boldsymbol{\mu}) \in X_{\mathcal{N}_t}^D$  such that

$$(6.17) \quad a(u_{\mathcal{N}_t}, v, \boldsymbol{\mu}) = 0, \quad \forall v \in X_{\mathcal{N}_t}.$$

### Output evaluation

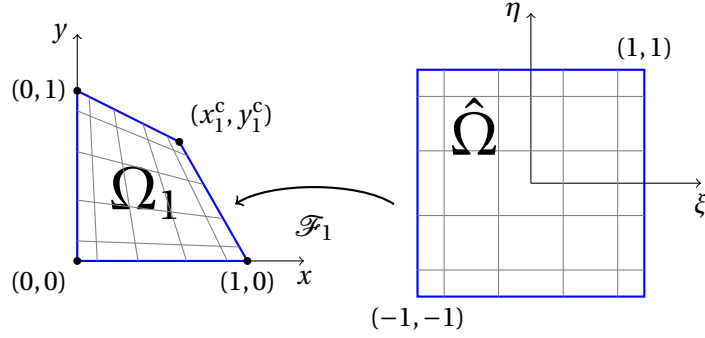
Our truth output of interest is essentially the six elements of the upper triangular part of the capacitance matrix, which implies solving (6.17) for  $\boldsymbol{\mu}$ ,  $\hat{\boldsymbol{\mu}}$  and  $\hat{\boldsymbol{\mu}}$ . For simplicity, we only consider evaluation of the first row of  $C$  – the outputs made available by  $u_{\mathcal{N}_t}(\boldsymbol{\mu})$  – in the discussion below.

Instead of calculating the outward normal derivative of  $u_{\mathcal{N}_t}(\boldsymbol{\mu})$  across  $\Gamma_1$ ,  $\Gamma_2$  and  $\Gamma_3$  in the integrals (6.4) and then integrate, we shall invoke the Neumann-Dirichlet equivalence described in Section 2.5 to evaluate our output of interest. We recall that with expanded spaces

$$(6.18) \quad \tilde{X}_{\mathcal{N}_t, j} \supset X_{\mathcal{N}_t}, \quad j = 1, 2, 3,$$

and under the assumptions of Lemma 2.1, we may formulate Neumann problems that are equivalent to the Dirichlet problems we actually solve and then profitably evaluate the flux integrals through the bilinear form. It is now necessary to formulate a different Neumann problem for each flux integral, hence the subscript  $j$  in the expansions above. As approximations to the (nonzero) elements in the first row of  $C(\boldsymbol{\mu})$ , we define

$$(6.19) \quad (C(\boldsymbol{\mu}))_{0, j} \approx l_j^{\text{out}}(u_{\mathcal{N}_t}(\boldsymbol{\mu})) \stackrel{\text{def}}{=} -a(u_{\mathcal{N}_t}(\boldsymbol{\mu}), v_j^*), \quad 1 \leq j \leq 3,$$



**Figure 6.4:** The mapping  $\mathcal{F}_1$  of  $\hat{\Omega}$  onto  $\Omega_1$ . The GLL-nodes, and the corresponding mapped nodes, are shown for the particular case  $P = 5$ .

and hence as a truth output of interest vector

$$(6.20) \quad \underline{s}_{\mathcal{N}_t}(\boldsymbol{\mu}) = \left[ l_1^{\text{out}}(u_{\mathcal{N}_t}(\boldsymbol{\mu})), l_2^{\text{out}}(u_{\mathcal{N}_t}(\boldsymbol{\mu})), l_3^{\text{out}}(u_{\mathcal{N}_t}(\boldsymbol{\mu})) \right]^T.$$

From Section 2.5, we recall that the functions  $v_j^*$  are equal to unity on  $\Gamma_j$  and identically zero on  $\Gamma_i$ ,  $i \neq j$ . Moreover, the particular choices of  $v_j^*$  implicitly define the expansions (6.18).

From Lemma 2.2, we see that as long as  $v_j^*$ ,  $j = 1, 2, 3$ , are chosen as members of  $V(\Omega)$  as defined in (6.16), every choice is equivalent. In fact, as GLL quadrature is used to numerically evaluate the integrals, every function will be “seen” as a polynomial of degree  $P_t$  from the numerical method’s point of view. Hence, in the spectral element context, every choice of  $v_j^*$  is discretely equivalent.

### Representation on the reference domain

As earlier mentioned (see Section 2.3), the computational realisation of the spectral element method involves the representation of the bilinear form  $a$  and the linear functional  $l$  on the reference domain  $\hat{\Omega}$  via the mappings  $\mathcal{F}_k^{-1}$ . As an example, Figure 6.4 illustrates the mapping  $\mathcal{F}_1 : \hat{\Omega} \rightarrow \Omega_1$ . A natural choice of mapping function is here the linear function

$$(6.21) \quad \mathcal{F}_1 : \begin{cases} x_1(\xi, \eta) = \frac{\xi+1}{2} \left( 1 + (x_1^c - 1) \frac{\eta+1}{2} \right) \\ y_1(\xi, \eta) = \frac{\eta+1}{2} \left( 1 + (y_1^c - 1) \frac{\xi+1}{2} \right), \end{cases}$$

where  $(x_1^c, y_1^c) = (\bar{x} - 0.4, \bar{y} - 0.4)$  denotes the coordinates of the upper right corner of the element. Similar mappings are readily constructed for the remaining elements.

Albeit the possibility of explicitly constructing the mapping functions, we use the automatic transfinite mapping algorithm proposed by Gordon and Hall [6] in our actual implementation. Thus, the code is readily extendable to more advanced geometrical shapes. However, as the boundaries of each of our subdomains are piecewise linear, the mapping constructed by Gordon-Hall is equivalent to the one above (for element number one).

Let us now explicitly consider the representation of the bilinear form on the reference domain. Clearly, we may write

$$(6.22) \quad a(u, v; \boldsymbol{\mu}) = \sum_{k=1}^9 \epsilon_k \int_{\Omega_k} \nabla u \cdot \nabla v \, d\Omega_k.$$

where  $\epsilon_k$  is equal to  $\epsilon_{\text{ano}}$  for  $k = 5$  and equal to  $\epsilon_{\text{bg}}$  otherwise. Considering term number  $k$ , we may after some algebra write

$$(6.23) \quad \int_{\Omega_1} \nabla u \cdot \nabla v \, d\Omega = \int_{\hat{\Omega}} (\hat{\nabla} \hat{u}_k)^T G_k(\boldsymbol{\mu}) \hat{\nabla} \hat{v}_k \, d\hat{\Omega},$$

where  $\hat{\nabla} \stackrel{\text{def}}{=} \left( \frac{\partial}{\partial \xi}, \frac{\partial}{\partial \eta} \right)^T$  and  $G_k \in \mathbb{R}^{2 \times 2}$  is determined by  $\mathcal{F}_k$  as

$$(6.24) \quad \begin{aligned} G_k &= \frac{1}{\det(J_k)} \begin{pmatrix} \tilde{g}_{k,11} & \tilde{g}_{k,12} \\ \tilde{g}_{k,21} & \tilde{g}_{k,22} \end{pmatrix} \\ &= \frac{1}{\det(J_k)} \begin{pmatrix} \left( \frac{\partial y_k}{\partial \eta} \right)^2 + \left( \frac{\partial x_k}{\partial \eta} \right)^2 & -\frac{\partial y_k}{\partial \xi} \frac{\partial y_k}{\partial \eta} - \frac{\partial x_k}{\partial \eta} \frac{\partial x_k}{\partial \xi} \\ -\frac{\partial y_k}{\partial \xi} \frac{\partial y_k}{\partial \eta} - \frac{\partial x_k}{\partial \eta} \frac{\partial x_k}{\partial \xi} & \left( \frac{\partial y_k}{\partial \xi} \right)^2 + \left( \frac{\partial x_k}{\partial \xi} \right)^2 \end{pmatrix} \end{aligned}$$

where  $\det(J_k) = \frac{\partial x_k}{\partial \xi} \frac{\partial y_k}{\partial \eta} - \frac{\partial x_k}{\partial \eta} \frac{\partial y_k}{\partial \xi}$  is the determinant of the Jacobian matrices.

Unfortunately, the functions  $g_{k,ij}(\boldsymbol{\mu})$  are for  $k = 1, 3, 7, 9$  non-affine. Hence, a rapid reduced basis offline-online computational decoupling will not be possible unless we may successfully employ empirical interpolation to our problem.

### Discrete reciprocity

Let us depart somewhat from our main path and further discuss the reciprocity of the physical system. Although being clear from the definition of  $C$ , the reciprocity may also be viewed as a direct consequence of the symmetry of the bilinear form.

Let  $u_0(\boldsymbol{\mu})$  correspond to the (exact) solution with a unity potential on electrode 0 for a given  $\boldsymbol{\mu}$ , and let  $u_1(\boldsymbol{\mu})$  correspond to the solution with a unity potential

## 6. A worked example: Electrostatics

---

on electrode 1 for the same parameter vector. From the definition of  $C$ , we have  $(C(\boldsymbol{\mu}))_{0,1} = (C(\boldsymbol{\mu}))_{1,0}$ , and since  $u_0$  and  $u_1$  are exact solutions, we have  $(C(\boldsymbol{\mu}))_{0,1} = -a(u_0, v_1^*; \boldsymbol{\mu})$  and  $(C(\boldsymbol{\mu}))_{1,0} = -a(u_1, v_0^*; \boldsymbol{\mu})$ . But now, as  $u_0$  and  $u_1$  belong to  $H^1$  and are equal to unity on  $\Gamma_0$  and  $\Gamma_1$ , respectively, we may set  $v_0^* = u_0$  and  $v_1^* = u_1$  and write

$$(6.25) \quad \begin{aligned} (C(\boldsymbol{\mu}))_{0,1} &= -a(u_0, v_1^*; \boldsymbol{\mu}) = -a(u_0, u_1; \boldsymbol{\mu}) \\ &= -a(u_1, u_0; \boldsymbol{\mu}) = -a(u_1, v_0^*; \boldsymbol{\mu}) = (C(\boldsymbol{\mu}))_{1,0}. \end{aligned}$$

Note that the first and last equalities are justified as we are working with exact solutions.

For spectral element approximations  $u_{\mathcal{N},0} \approx u_0$  and  $u_{\mathcal{N},1} \approx u_1$ , the above result still holds, as we may indeed choose  $v_0^* = u_{\mathcal{N},0}$  and  $v_1^* = u_{\mathcal{N},1}$  by Lemma 2.2. Hence, the spectral element (and any other standard finite element) approximation preserves the reciprocity of the physical system.

### 6.2.3 RB formulation

Armed with the RB framework from Chapter 4, the definition (6.6) of the parameter space  $\mathcal{D}$  and the parametric weak form (6.13), we are ready to formulate the reduced basis problem. We shall formulate our problem homogeneously, as this will prove easier to work with within the RB context. To this end, we define a boundary lifting function  $u^D(\boldsymbol{\mu}) \in X_{\mathcal{N}_t}^D$ , which is equal to 1 on  $\Gamma_0$ , and equal to zero at every interior node, and write  $u_{\mathcal{N}_t}(\boldsymbol{\mu}) = u_{\mathcal{N}_t}^D + u_{\mathcal{N}_t}^0(\boldsymbol{\mu})$  with  $u_{\mathcal{N}_t}^0 \in X_{\mathcal{N}_t}$ . Given a set of parameter vectors  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_{N_{\max}}$ , our precomputed snapshots are now given as  $u_{\mathcal{N}_t}^0(\boldsymbol{\mu}_n)$ , for  $1 \leq n \leq N_{\max}$ .

As our RB approximation spaces, we define

$$(6.26) \quad X_N \stackrel{\text{def}}{=} \text{span} \left\{ u_{\mathcal{N}_t}^0(\boldsymbol{\mu}_n) \right\}_{n=1}^N, \quad 1 \leq N \leq N_{\max}.$$

Given a parameter vector  $\boldsymbol{\mu} \in (\mathcal{D} \cup \overline{\mathcal{D}})$ , we shall require the RB solutions

$$(6.27) \quad u_N(\boldsymbol{\kappa}) \text{ or } u_N(\overline{\boldsymbol{\kappa}}) \quad \text{for all } \boldsymbol{\kappa} \in (\boldsymbol{\mu}, \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\mu}}),$$

to evaluate our RB output *matrix* of interest. We choose  $\boldsymbol{\kappa}$  or  $\overline{\boldsymbol{\kappa}}$ , dependent upon whether it is  $\boldsymbol{\kappa}$  or  $\overline{\boldsymbol{\kappa}}$  that belongs to  $\mathcal{D}$  (we recall  $\overline{\boldsymbol{\kappa}}$  as the “flipping” of  $\boldsymbol{\kappa}$  across the vertical centreline of  $\Omega$ ).

The reduced problem reads as follows: Given any  $\boldsymbol{\mu} \in (\mathcal{D} \cup \overline{\mathcal{D}})$ , then for  $\boldsymbol{\kappa} \in (\boldsymbol{\mu}, \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\mu}})$ , check whether  $\boldsymbol{\kappa} \in \mathcal{D}$ . If it is, find  $u_N^0(\boldsymbol{\kappa}) = u_N(\boldsymbol{\kappa}) - u^D(\boldsymbol{\kappa})$  such that

$$(6.28) \quad a(u_N, v; \boldsymbol{\kappa}) = 0, \quad \forall v \in X_N.$$

If on the other hand  $\boldsymbol{\kappa} \in \overline{\mathcal{D}}$ , find  $u_N^0(\overline{\boldsymbol{\kappa}}) = u_N(\overline{\boldsymbol{\kappa}}) - u^D(\overline{\boldsymbol{\kappa}})$  such that

$$(6.29) \quad a(u_N, v; \overline{\boldsymbol{\kappa}}) = 0, \quad \forall v \in X_N.$$

Finally, evaluate the RB output of interest *matrix*, defined as

$$(6.30) \quad s_N(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \begin{bmatrix} l_1^{\text{out}}(u_N(\boldsymbol{\mu})) & l_2^{\text{out}}(u_N(\boldsymbol{\mu})) & l_3^{\text{out}}(u_N(\boldsymbol{\mu})) \\ 0 & l_2^{\text{out}}(u_N(\hat{\boldsymbol{\mu}})) & l_3^{\text{out}}(u_N(\hat{\boldsymbol{\mu}})) \\ 0 & 0 & l_3^{\text{out}}(u_N(\hat{\boldsymbol{\mu}})) \end{bmatrix}.$$

We shall refer to  $u_N$  and  $s_N$  as the ‘‘RB solution’’ and ‘‘RB output’’, respectively.

#### 6.2.4 RB formulation with Empirical Interpolation (RB-EI)

We now invoke the EI method in the RB formulation, and start by recalling that the bilinear form can be written in terms of the reference variables as

$$(6.31) \quad a(u, v; \boldsymbol{\mu}) = \sum_{k=1}^9 \epsilon_k \int_{\hat{\Omega}} (\hat{\mathbf{V}} \hat{u}_k)^T G_k(\boldsymbol{\mu}) \hat{\mathbf{V}} \hat{v}_k \, d\hat{\Omega},$$

where

$$(6.32) \quad G_k(\boldsymbol{\mu}) = \begin{bmatrix} g_{k,11}(\boldsymbol{\mu}) & g_{k,12}(\boldsymbol{\mu}) \\ g_{k,21}(\boldsymbol{\mu}) & g_{k,22}(\boldsymbol{\mu}) \end{bmatrix}$$

are matrices of geometrical factors corresponding to the mappings  $\mathcal{F}_k$ . Here,  $g_{k,ij} = \tilde{g}_{k,ij} / \det(J_k)$ , where the  $\tilde{g}_{k,ij}$  are given in (6.24) and the  $\det(J_k)$  are the determinants of the Jacobian matrices. Expanding the first term of (6.31) (and setting  $\epsilon_1 = 1$ ), we get

$$(6.33) \quad \begin{aligned} & \int_{\hat{\Omega}} (\hat{\mathbf{V}} \hat{u}_1)^T G_1(\boldsymbol{\mu}) \hat{\mathbf{V}} \hat{v}_1 \, d\hat{\Omega} \\ &= \int_{\hat{\Omega}} \left( g_{1,11}(\boldsymbol{\mu}) \frac{\partial \hat{u}_1}{\partial \xi} \frac{\partial \hat{v}_1}{\partial \xi} + g_{1,12}(\boldsymbol{\mu}) \frac{\partial \hat{u}_1}{\partial \xi} \frac{\partial \hat{v}_1}{\partial \eta} \right. \\ & \quad \left. + g_{1,21}(\boldsymbol{\mu}) \frac{\partial \hat{u}_1}{\partial \eta} \frac{\partial \hat{v}_1}{\partial \xi} + g_{1,22}(\boldsymbol{\mu}) \frac{\partial \hat{u}_1}{\partial \eta} \frac{\partial \hat{v}_1}{\partial \eta} \right) d\hat{\Omega}. \end{aligned}$$

## 6. A worked example: Electrostatics

---

To recover an efficient offline-online computational decoupling, the non-affine functions  $g_{1,11}(\boldsymbol{\mu}), g_{1,12}(\boldsymbol{\mu}), g_{1,21}(\boldsymbol{\mu}), g_{1,22}(\boldsymbol{\mu})$  need to be affinely, and independently, approximated with the EI method.

As it turns out, only the functions  $g_{k,ij}$  for  $k = 1, 3, 7, 9$  (corresponding to the four corner elements) are non-affine. In fact, a total of 16 functions are in need of empirical interpolation, whereas the remaining terms may be written affinely in 36 terms. Hence, our approximate affine representation of the bilinear form has  $Q_M = 16M + 36$  terms, where  $M$  is the number of (empirical) interpolation nodes. Abstractly, we write

$$(6.34) \quad a(\cdot, \cdot; \boldsymbol{\mu}) \approx a_M(\cdot, \cdot; \boldsymbol{\mu}) \stackrel{\text{def}}{=} \sum_{q=1}^{Q_M} \Theta_{a_M}^q(\boldsymbol{\mu}) a_M^q(\cdot, \cdot).$$

The reduced ‘‘RB-EI’’ problem now reads: Given any  $\boldsymbol{\mu} \in (\mathcal{D} \cup \overline{\mathcal{D}})$ , then for each  $\boldsymbol{\kappa} \in (\boldsymbol{\mu}, \widehat{\boldsymbol{\mu}}, \widehat{\boldsymbol{\mu}})$ , check whether  $\boldsymbol{\kappa} \in \mathcal{D}$ . If it is, find  $u_N^{M,0}(\boldsymbol{\kappa}) = u_N^M(\boldsymbol{\kappa}) - u^D(\boldsymbol{\kappa}) \in X_N$  such that

$$(6.35) \quad a_M(u_N^M, v; \boldsymbol{\kappa}) = 0, \quad \forall v \in X_N.$$

If on the other hand  $\boldsymbol{\kappa} \in \overline{\mathcal{D}}$ , find  $u_N^{M,0}(\overline{\boldsymbol{\kappa}}) = u_N^M(\overline{\boldsymbol{\kappa}}) - u^D(\overline{\boldsymbol{\kappa}}) \in X_N$  such that

$$(6.36) \quad a_M(u_N^M, v; \overline{\boldsymbol{\kappa}}) = 0, \quad \forall v \in X_N.$$

Finally, there are now two obvious ways to evaluate the output of interest. Either, we evaluate the output matrix as

$$(6.37) \quad s_N^M(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \begin{bmatrix} l_1^{\text{out}}(u_N^M(\boldsymbol{\mu})) & l_2^{\text{out}}(u_N^M(\boldsymbol{\mu})) & l_3^{\text{out}}(u_N^M(\boldsymbol{\mu})) \\ 0 & l_2^{\text{out}}(u_N^M(\widehat{\boldsymbol{\mu}})) & l_3^{\text{out}}(u_N^M(\widehat{\boldsymbol{\mu}})) \\ 0 & 0 & l_3^{\text{out}}(u_N^M(\widehat{\boldsymbol{\mu}})) \end{bmatrix},$$

i.e., through the bilinear form  $a$ , or as

$$(6.38) \quad \tilde{s}_N^M(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \begin{bmatrix} l_{M,1}^{\text{out}}(u_N^M(\boldsymbol{\mu})) & l_{M,2}^{\text{out}}(u_N^M(\boldsymbol{\mu})) & l_{M,3}^{\text{out}}(u_N^M(\boldsymbol{\mu})) \\ 0 & l_{M,2}^{\text{out}}(u_N^M(\widehat{\boldsymbol{\mu}})) & l_{M,3}^{\text{out}}(u_N^M(\widehat{\boldsymbol{\mu}})) \\ 0 & 0 & l_{M,3}^{\text{out}}(u_N^M(\widehat{\boldsymbol{\mu}})) \end{bmatrix},$$

where

$$(6.39) \quad l_{M,j}^{\text{out}}(u(\boldsymbol{\mu})) \stackrel{\text{def}}{=} -a_M(u, v_j^*; \boldsymbol{\mu}), \quad 1 \leq j \leq 3,$$

i.e. through the bilinear form  $a_M$ . As the empirical interpolation error tends to zero, we expect  $l_{M,j}^{\text{out}}(u(\boldsymbol{\mu})) \rightarrow l_j^{\text{out}}(u(\boldsymbol{\mu}))$ .

On occasion, we shall refer to  $u_N^M$  and  $\tilde{s}_N^M$  as the “RB-EI solution” and “RB-EI output”, respectively.

Finally, we emphasise that as  $a_M(\cdot, \cdot; \boldsymbol{\mu})$  is an affine bilinear form, our problem now admits an efficient offline-online computational decoupling, as described in Section 4.2.

### 6.2.5 Remarks on RB and RB-EI output evaluation

Below, we again only consider the outputs corresponding to the “unrotated” parameter vector  $\boldsymbol{\mu}$ , now given as the RB output vector

$$(6.40) \quad \underline{s}_N(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \left[ l_1^{\text{out}}(u_N(\boldsymbol{\mu})), l_2^{\text{out}}(u_N(\boldsymbol{\mu})), l_3^{\text{out}}(u_N(\boldsymbol{\mu})) \right]^T.$$

We first discuss the choice of the functions  $v_j^*$ , which play important roles in the evaluation of both the RB and RB-EI outputs. Our discussion deals with the RB-output, but applies to the RB-EI output as well. Then, we describe an efficient way to evaluate the RB-EI output.

#### Choice of $v^*$

When evaluating the  $l_j^{\text{out}}$  of (6.40), two arbitrary choices  $v_{j,1}^*, v_{j,2}^*$  for  $v_j^*$  (recall the definitions (6.19) of the output functionals) are now not equivalent as was the case for the SE approximation. In fact, this is a consequence of Lemma 2.2, as now  $w^* = v_{j,1}^* - v_{j,2}^*$  does not in general belong to  $X_N$ .

As our outputs of interest are evaluated through the bilinear form, the error estimate (2.35), restated here as

$$(6.41) \quad \left| l_j^{\text{out}}(u_N(\boldsymbol{\mu})) - l_j^{\text{out}}(u_{\mathcal{N}_t}(\boldsymbol{\mu})) \right| \leq (a(e_N, e_N; \boldsymbol{\mu}))^{1/2} (a(v_j^*, v_j^*; \boldsymbol{\mu}))^{1/2},$$

for  $j = 1, 2, 3$ , suggests choosing  $v_j^*$  such that

$$(6.42) \quad v_j^* = \arg \min_{v \in V_j^*} a(v, v; \boldsymbol{\mu})$$

where

$$(6.43) \quad V_j^* \stackrel{\text{def}}{=} \{v \in V : v|_{\{\Gamma_i\}_{i=0}^3 \setminus \Gamma_j} = 0, v|_{\Gamma_j} = 1\}, \quad j = 1, 2, 3,$$

## 6. A worked example: Electrostatics

---

and  $V$  is the truth approximation space defined in (6.16). As  $a$  is symmetric,  $v_j^* \in V_j^*$  is the solution to [25]

$$(6.44) \quad a(v_j^*, v; \boldsymbol{\mu}) = 0, \quad \forall v \in (V \cap \{v : v|_{\{\Gamma_i\}_{i=1}^3} = 0\}).$$

Obviously, we cannot solve a (three, in fact) problems of “truth complexity” at every evaluation of the RB output matrix, and thus good surrogates for the  $v_j^*$  need to be found. We shall consider three alternatives:

- i)* Compute  $v_j^*$  as the solution to (6.44) using approximation spaces of (very) low order, e.g. quadratic or cubic polynomials.
- ii)* Compute  $v_j^*$  as the solution to (6.44) with  $\boldsymbol{\mu}$  replaced by  $\boldsymbol{\mu}_{\text{ref}}$  – thus the  $v_j^*$  can be computed as part of the preprocessing stage.
- iii)* Simply choose  $v_j^*$  equal to unity on  $\Gamma_j$  and equal to zero in all other nodes. We note that this choice will be far from the minimiser (6.42), as  $v_j^*$  comprise a large amount of “energy” in the vicinity of  $\Gamma_j$ .

### Efficient output evaluation

To compute the elements of the output matrix  $s_N$  defined in (6.30),  $s_N^M$  defined in (6.37) or  $\tilde{s}_N^M$  defined in (6.38), the straightforward approach would be constructing the RB or RB-EI solution and the  $v_j^*$  explicitly and then evaluate the six required inner-products. Such a procedure however, is  $\mathcal{N}_t$ -dependent and thus encourages the pursuit of a more efficient alternative.

Our aim is to develop an entirely  $\mathcal{N}_t$ -independent online procedure for output evaluation. Hence, it seems reasonable to put the matrices  $s_N$  and  $s_N^M$  aside as they comprise nonaffine  $a$ -inner-products and instead work with  $\tilde{s}_N^M$ , in which the output functionals are evaluated through the affine bilinear form  $a_M$ .

For the sake of simplifying our discussion, we now only consider evaluation of the single, scalar-valued output functional  $l_{M,1}^{\text{out}}$  defined in (6.39), and assume that  $\boldsymbol{\mu}$  resides in  $\mathcal{D}$ . The algebraic formulation of the RB-EI problem may then be stated as: Find  $\underline{u}_N^{M,0}(\boldsymbol{\mu}) = \underline{u}_N^M(\boldsymbol{\mu}) - u^D(\boldsymbol{\mu}) \in X_N$  subject to

$$(6.45) \quad \sum_{n=1}^N u_{N,n}^{M,0}(\boldsymbol{\mu}) \underbrace{a_M(\zeta^n, \zeta^m; \boldsymbol{\mu})}_{(A_N^M(\boldsymbol{\mu}))_{mn}} = \underbrace{-a_M(u^D, \zeta^m; \boldsymbol{\mu})}_{(l_N^M(\boldsymbol{\mu}))_m}, \quad 1 \leq m \leq N,$$

where  $u_N^{M,0}(x, y; \boldsymbol{\mu}) = \sum_{n=1}^N u_{N,n}^{M,0}(\boldsymbol{\mu}) \zeta^n(x, y)$ . Equivalently, we may write (6.45) as

$$(6.46) \quad A_N^M(\boldsymbol{\mu}) \underline{u}_N^{M,0}(\boldsymbol{\mu}) = \underline{l}_N^M(\boldsymbol{\mu}),$$



where

$$(6.47) \quad \underline{u}_N^{M,0}(\boldsymbol{\mu}) = [u_{N,1}^{M,0}(\boldsymbol{\mu}), \dots, u_{N,N}^{M,0}(\boldsymbol{\mu})]^\top \in \mathbb{R}^N$$

is the RB-EI vector of unknowns (and the elements of  $A_N^M$  and  $\underline{l}_N^M$  are indicated in (6.45)).

We now define an augmented RB-EI solution vector

$$(6.48) \quad \tilde{\underline{u}}_N^M(\boldsymbol{\mu}) \stackrel{\text{def}}{=} [1, 0, u_{N,1}^{M,0}(\boldsymbol{\mu}), \dots, u_{N,N}^{M,0}(\boldsymbol{\mu})]^\top \in \mathbb{R}^{N+2},$$

and write the RB-EI solution as

$$(6.49) \quad \begin{aligned} u_N^M(x, y; \boldsymbol{\mu}) &= u_N^{M,0}(x, y; \boldsymbol{\mu}) + u^D(x, y; \boldsymbol{\mu}) \\ &= \sum_{n=-1}^N (\tilde{\underline{u}}_{N,n}^M(\boldsymbol{\mu}))_{n+2} \zeta^n(x, y), \end{aligned}$$

where  $\zeta^{-1} \stackrel{\text{def}}{=} u^D$  and  $\zeta^0 \stackrel{\text{def}}{=} v_1^*$  (note that the term for  $n = 0$  always vanishes). We also define an augmented RB-EI stiffness matrix  $\tilde{A}_N^M \in \mathbb{R}^{N+2 \times N+2}$  given by

$$(6.50) \quad \tilde{A}_N^M(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \begin{bmatrix} a_M(u^D, u^D; \boldsymbol{\mu}) & a_M(v_1^*, u^D; \boldsymbol{\mu}) & a_M(\zeta^1, u^D; \boldsymbol{\mu}) & \dots & a_M(\zeta^N, u^D; \boldsymbol{\mu}) \\ a_M(u^D, v_1^*; \boldsymbol{\mu}) & a_M(v_1^*, v_1^*; \boldsymbol{\mu}) & a_M(\zeta^1, v_1^*; \boldsymbol{\mu}) & \dots & a_M(\zeta^N, v_1^*; \boldsymbol{\mu}) \\ a_M(u^D, \zeta^1; \boldsymbol{\mu}) & a_M(v_1^*, \zeta^1; \boldsymbol{\mu}) & a_M(\zeta^1, \zeta^1; \boldsymbol{\mu}) & \dots & a_M(\zeta^1, \zeta^1; \boldsymbol{\mu}) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_M(u^D, \zeta^N; \boldsymbol{\mu}) & a_M(v_1^*, \zeta^N; \boldsymbol{\mu}) & a_M(\zeta^1, \zeta^N; \boldsymbol{\mu}) & \dots & a_M(\zeta^N, \zeta^N; \boldsymbol{\mu}) \end{bmatrix},$$

where we in the two first rows and columns have included the terms resulting from incorporating  $\zeta^{-1} = u^D$  and  $\zeta^0 = v_1^*$  in the basis for the RB-EI solution.

Indeed, we may construct  $\tilde{A}_N^M$  following the offline-online procedure described for affine problems in Section 4.2. Analogously to (4.35), we now invoke (6.34) and write

$$(6.51) \quad \begin{aligned} (\tilde{A}_N^M(\boldsymbol{\mu}))_{m+2, n+2} &= \sum_{i=1}^{\mathcal{N}_t} \sum_{j=1}^{\mathcal{N}_t} \zeta_j^m \zeta_i^n a_M(\psi_j, \psi_i; \boldsymbol{\mu}) \\ &= \sum_{i=1}^{\mathcal{N}_t} \sum_{j=1}^{\mathcal{N}_t} \zeta_j^m \zeta_i^n \left( \sum_{q=1}^{Q_M} \Theta_{a_M}^q(\boldsymbol{\mu}) a_M^q(\psi_j, \psi_i) \right) \\ &= \sum_{q=1}^{Q_M} \Theta_{a_M}^q(\boldsymbol{\mu}) \underbrace{\sum_{i=1}^{\mathcal{N}_t} \sum_{j=1}^{\mathcal{N}_t} \zeta_j^m \zeta_i^n a_M^q(\psi_j, \psi_i)}_{=(\tilde{A}_N^{M,q})_{m+2, n+2}}, \quad -1 \leq m, n \leq N, \end{aligned}$$

## 6. A worked example: Electrostatics

---

where the  $N + 2 \times N + 2$  matrices  $\tilde{A}_N^{M,q}$  are parameter independent and pre-computable. Note that the RB-EI stiffness matrix,  $A_N^M(\boldsymbol{\mu})$ , is simply the  $N \times N$  lower right submatrix of  $\tilde{A}_N^M$  and that  $a_M$  – and consequently  $\tilde{A}_N^M$  – is symmetric.

Finally, computing the “residual vector”  $\underline{r}_N^M = \tilde{A}_N^M \underline{\tilde{u}}_N^M$  (compare to the original equation (6.46), reading  $A_N^M \underline{u}_N^{M,0} - \underline{l}_N^M = 0$ ), and then taking

$$(6.52) \quad l_{M,1}^{\text{out}} = -(\underline{r}_N^M)_2 \left( = -a_M(u_N^M, v_1^*; \boldsymbol{\mu}) \right),$$

is presumably much quicker than direct evaluation of  $a_M(u_N^M, v_1^*; \boldsymbol{\mu})$  once the parameter independent matrices  $\tilde{A}_N^{M,q}$  have been formed. In fact, we assemble  $\tilde{A}_N^M$  (and thus also  $A_N^M$ ) in  $\mathcal{O}(Q_M N^2)$  flops. Note that in actual practice, we may compute the element  $(\underline{r}_N^M)_2$  directly as the sum

$$(6.53) \quad (\underline{r}_N^M)_2 = \sum_{j=1}^{N+2} (\tilde{A}_N^M)_{2,j} (\tilde{\underline{u}}_N^M)_j$$

in  $\mathcal{O}(N + 2)$  operations (i.e. without performing the full  $(N + 2)^2$ -flops operator evaluation).

### 6.2.6 Remarks on inhomogeneous Dirichlet conditions

Inhomogeneous Dirichlet boundary conditions are seldomly considered in the existing RB literature. Although quite straightforward in the standard FE or SE methods, the imposition of such boundary conditions turns out rather peculiar in the RB context.

Above, we made an “obvious” choice of boundary lifting function  $u^D$ , simply by choosing  $u^D$  equal to unity on  $\Gamma_0$  and equal to zero at every interior node. For finite- or spectral element methods, every choice of boundary lifting is equivalent due to the richness of the approximation spaces. In the RB context however, every choice is not equivalent but does in fact define a new homogeneous problem. To render this more apparent, let

$$(6.54) \quad u_{\mathcal{N}_t} = u^{\text{D},1} + u_{\mathcal{N}_t}^{0,1} = u^{\text{D},2} + u_{\mathcal{N}_t}^{0,2},$$

where  $u^{\text{D},1}, u^{\text{D},2} \in X_{\mathcal{N}_t}^{\text{D}}$  are different liftings of the Dirichlet data. For a set of parameters  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_N$ , the corresponding RB spaces are

$$(6.55) \quad X_N^1 = \text{span}\{u^{0,1}(\boldsymbol{\mu}_n)\}_{n=1}^N, \quad X_N^2 = \text{span}\{u^{0,2}(\boldsymbol{\mu}_n)\}_{n=1}^N.$$

In other words, the different choices of lifting give rise to two RB problems: Given  $\boldsymbol{\mu} \in \mathcal{D}$ , find  $u_N^{0,i}(\boldsymbol{\mu}) = u_N^i(\boldsymbol{\mu}) - u^{\text{D},i}(\boldsymbol{\mu}) \in X_N^i$  such that

$$(6.56) \quad a(u_N^i, v; \boldsymbol{\mu}) = 0, \quad \forall v \in X_N^i,$$

where  $i = 1, 2$ . As it turns out then, our choice was not obvious at all! However, every choice of  $u^{\text{D}} \in X_{\mathcal{N}_t}^{\text{D}}$  will be valid in the sense that the corresponding RB problem is well defined, and that the RB space is tailored to the specific problem at hand.

We leave this subject merely as a curious notice, and henceforth hold on to our choice of  $u^{\text{D}}$ .

### 6.2.7 *A posteriori* error estimation

As regards *a posteriori* error estimation, we confine ourselves within the framework for affine problems described in Section 4.4. As a consequence, our estimators *i)* do not rigorously apply when empirical interpolation is used, and *ii)* require the solution of a truth-complexity problem at each evaluation, and are thus practically useless.

However, following the procedure in [15, Chapter 6], constructing *a posteriori* estimators for the RB solution  $u_N^M$  (using empirical interpolation) seems to be within reach both theoretically and practically. Moreover, were we to construct such estimators, we might plausibly assume them to resemble the behavior of the standard estimators as the interpolation error tends to zero. Theoretically speaking then, developing (and later on numerically examining) the standard estimators may provide useful insight.

#### A lower bound for the coercivity constant

We now develop a lower bound  $\alpha_{\text{LB}}$  for the coercivity constant  $\alpha$ . Firstly, our constant must comply to

$$(6.57) \quad \alpha_{\text{LB}}(\boldsymbol{\mu}) \leq \alpha(\boldsymbol{\mu}) = \inf_{v \in X_{\mathcal{N}_t}} \frac{a(v, v; \boldsymbol{\mu})}{a(v, v; \boldsymbol{\mu}_{\text{ref}})},$$

and secondly, for our *a posteriori* error estimator to be sharp,  $\alpha_{\text{LB}}$  should not be a too conservative lower bound. We note that  $\alpha_{\text{LB}}$  is a factor both in the standard estimators (Chapter 4) and in the estimator that also incorporate the error resulting from empirical interpolation [15].

## 6. A worked example: Electrostatics

---

Below, we make use of a procedure similar to one in [13], although in a slightly different context. We recall that we may write the bilinear form  $a(u, v, \boldsymbol{\mu})$  in terms of the reference variables as

$$(6.58) \quad a(u, v; \boldsymbol{\mu}) = \sum_{k=1}^9 \epsilon_k \int_{\hat{\Omega}} (\hat{\mathbf{V}} \hat{u}_k)^T G_k(\boldsymbol{\mu}) \hat{\mathbf{V}} \hat{v}_k \, d\hat{\Omega},$$

where the  $G_k(\boldsymbol{\mu})$  are parameter dependent  $2 \times 2$  matrices comprising geometrical factors. From (6.24), we know that the  $G_k(\boldsymbol{\mu})$  are symmetric, so we may write  $G_k = Q_k^T \Lambda_k Q_k$  where the  $\Lambda_k$  are diagonal matrices with the eigenvalues  $\lambda_k^{\max} \geq \lambda_k^{\min}$  of  $G_k$  as elements, and the  $Q_k$  are orthogonal matrices (i.e.  $Q_k^T Q_k = I$  is the identity matrix) with the eigenvectors of  $G_k$  as column vectors. For  $v, w \in X$ , we may now write

$$(6.59) \quad a(v, w; \boldsymbol{\mu}) = \sum_{k=1}^9 \epsilon_k \int_{\hat{\Omega}} (\hat{\mathbf{V}} \hat{v}_k)^T G_k(\boldsymbol{\mu}) \hat{\mathbf{V}} \hat{w}_k \, d\hat{\Omega}$$

$$(6.60) \quad \geq \sum_{k=1}^9 \epsilon_k \int_{\hat{\Omega}} \lambda_k^{\min}(\boldsymbol{\mu}) (Q_k \hat{\mathbf{V}} \hat{v}_k)^T Q_k \hat{\mathbf{V}} \hat{w}_k \, d\hat{\Omega}$$

$$(6.61) \quad = \sum_{k=1}^9 \epsilon_k \int_{\hat{\Omega}} \lambda_k^{\min}(\boldsymbol{\mu}) (\hat{\mathbf{V}} \hat{v}_k)^T \hat{\mathbf{V}} \hat{w}_k \, d\hat{\Omega},$$

and similarly

$$(6.62) \quad a(v, w; \boldsymbol{\mu}) = \sum_{k=1}^9 \epsilon_k \int_{\hat{\Omega}} (\hat{\mathbf{V}} \hat{v}_k)^T G_k(\boldsymbol{\mu}) \hat{\mathbf{V}} \hat{w}_k \, d\hat{\Omega}$$

$$(6.63) \quad \leq \sum_{k=1}^9 \epsilon_k \int_{\hat{\Omega}} \lambda_k^{\max}(\boldsymbol{\mu}) (\hat{\mathbf{V}} \hat{v}_k)^T \hat{\mathbf{V}} \hat{w}_k \, d\hat{\Omega}.$$

With

$$(6.64) \quad \lambda^-(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \min_{\substack{(\xi, \eta) \in \hat{\Omega} \\ 1 \leq k \leq 9}} \lambda_k^{\min}(\xi, \eta; \boldsymbol{\mu})$$

and

$$(6.65) \quad \lambda^+(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \max_{\substack{(\xi, \eta) \in \hat{\Omega} \\ 1 \leq k \leq 9}} \lambda_k^{\max}(\xi, \eta; \boldsymbol{\mu}),$$

we may now for any  $v \in X$  and any two parameter vectors  $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2 \in \mathcal{D}$  write

$$(6.66) \quad \frac{a(v, v; \boldsymbol{\mu}_1)}{a(v, v; \boldsymbol{\mu}_2)} \geq \frac{\lambda^-(\boldsymbol{\mu}_1) \sum_{k=1}^9 \epsilon_k \int_{\hat{\Omega}} (\hat{\mathbf{V}} \hat{v}_k)^T \hat{\mathbf{V}} \hat{v}_k \, d\hat{\Omega}}{\lambda^+(\boldsymbol{\mu}_2) \sum_{k=1}^9 \epsilon_k \int_{\hat{\Omega}} (\hat{\mathbf{V}} \hat{v}_k)^T \hat{\mathbf{V}} \hat{v}_k \, d\hat{\Omega}} \\ = \frac{\lambda^-(\boldsymbol{\mu}_1)}{\lambda^+(\boldsymbol{\mu}_2)},$$

and hence in particular as a lower bound for the coercivity constant

$$(6.67) \quad \alpha_{\text{LB}} \stackrel{\text{def}}{=} \frac{\lambda^-(\boldsymbol{\mu})}{\lambda^+(\boldsymbol{\mu}_{\text{ref}})}.$$

### Error estimators

We first recall from Section 4.4 the general estimate for the energy-norm error of the field variable

$$(6.68) \quad \|e_N(\boldsymbol{\mu})\|_{\mathcal{E}} \leq \Delta_N(\boldsymbol{\mu}),$$

where  $e_N(\boldsymbol{\mu}) = u_N(\boldsymbol{\mu}) - u_{\mathcal{N}_t}(\boldsymbol{\mu})$  and, from Theorem 4.2,

$$(6.69) \quad \Delta_N(\boldsymbol{\mu}) = \frac{\|\hat{e}(\boldsymbol{\mu})\|_{\boldsymbol{\mu}_{\text{ref}}}}{\sqrt{\alpha_{\text{LB}}(\boldsymbol{\mu})}},$$

where  $\alpha_{\text{LB}}$  is the coercivity lower bound developed above. By invoking the output error estimate (2.35), we have, for  $1 \leq j \leq 3$ ,

$$(6.70) \quad \begin{aligned} \left| (\underline{s}_N(\boldsymbol{\mu}))_j - (\underline{s}_{\mathcal{N}_t}(\boldsymbol{\mu}))_j \right| &\leq \|e_N(\boldsymbol{\mu})\|_{\mathcal{E}} a(v_j^*, v_j^*; \boldsymbol{\mu}) \\ &\leq \Delta_N(\boldsymbol{\mu}) a(v_j^*, v_j^*; \boldsymbol{\mu}). \end{aligned}$$

As output error bounds, we thus define

$$(6.71) \quad \Delta_{N,j}^{\text{out}}(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \Delta_N(\boldsymbol{\mu}) a(v_j^*, v_j^*; \boldsymbol{\mu}), \quad 1 \leq j \leq 3.$$

As a maximum output error upper bound, we also introduce

$$(6.72) \quad \Delta_N^{\text{out,max}}(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \max_{1 \leq j \leq 3} \Delta_{N,j}^{\text{out}}(\boldsymbol{\mu}).$$

These estimators may now be used both as a bound for the error during the parameter vector selection, and for certification of the accuracy of the RB solution or outputs, albeit at “truth” computational cost.

## 6.3 An inverse problem

In the previous section, we considered the forward problem  $\boldsymbol{\mu} \rightarrow s(\boldsymbol{\mu})$ . In this section, the forward model is employed in solving the *inverse* problem of estimating the parameter vector that corresponds to a given matrix of output data.

## 6. A worked example: Electrostatics

---

We assume that we are given an upper triangular matrix  $s_{\text{obs}} \in \mathbb{R}^{3 \times 3}$  of capacitance measurements. Then, the inverse problem reads: Find  $\boldsymbol{\mu} \in (\mathcal{D} \cup \overline{\mathcal{D}})$  such that

$$(6.73) \quad \boldsymbol{\mu} = \arg \min_{\boldsymbol{\mu} \in \mathcal{D} \cup \overline{\mathcal{D}}} \|s_{\text{obs}} - s(\boldsymbol{\mu})\|,$$

for some norm  $\|\cdot\|$ . In general,  $s_{\text{obs}} = s(\boldsymbol{\mu})$  is unlikely due to *i*) noisy capacitance measurements, *ii*) numerical error in the forward model or *iii*) biased forward outputs due to e.g. wrong assumptions made in the mathematical model.

As our forward model, we choose  $\boldsymbol{\mu} \rightarrow \tilde{s}_N^M(\boldsymbol{\mu})$  and thus anticipate significant numerical errors for small  $M$  and  $N$ . We then define the residual

$$(6.74) \quad R_N^M(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \frac{1}{2} \sum_{i=1}^3 \sum_{j=1}^3 \left( (s_{\text{obs}})_{ij} - (\tilde{s}_N^M(\boldsymbol{\mu}))_{ij} \right)^2.$$

and seek the reconstructed anomaly position  $\boldsymbol{\mu}_{N,\text{rec}}^M$  as

$$(6.75) \quad \boldsymbol{\mu}_{N,\text{rec}}^M = \arg \min_{\boldsymbol{\mu} \in \mathcal{D} \cup \overline{\mathcal{D}}} R_N^M(\boldsymbol{\mu}).$$

We shall assume that  $R_N^M(\boldsymbol{\mu})$  attains a single minimum for  $\boldsymbol{\mu} \in (\mathcal{D} \cup \overline{\mathcal{D}})$ , that  $R_N^M$  has no saddle points and that the Jacobian of  $\tilde{s}_N^M(\boldsymbol{\mu})$  is invertible.

Given an appropriate initial value  $\boldsymbol{\mu}^0$ , the minimiser  $\boldsymbol{\mu}_N^M$  of the nonlinear least squares problem (6.75) is sought numerically by applying a Gauss-Newton iterative scheme [17]

$$(6.76) \quad \boldsymbol{\mu}^{k+1} = \boldsymbol{\mu}^k - (H_{R_N^M}(\boldsymbol{\mu}^k))^{-1} \nabla R_N^M(\boldsymbol{\mu}^k),$$

where  $H_{R_N^M}(\boldsymbol{\mu}^k) \in \mathbb{R}^{2 \times 2}$  and  $\nabla R_N^M(\boldsymbol{\mu}^k) \in \mathbb{R}^2$  denote the approximate Hessian and gradient of  $R_N^M$ , respectively, at the  $k$ 'th iterate  $\boldsymbol{\mu}^k$ . To approximate the derivatives in  $H_{R_N^M}$  and  $\nabla R_N^M$ , we apply the standard central difference formula

$$(6.77) \quad \frac{\partial f}{\partial x_i} \approx \frac{f(x_1, \dots, x_i + h, \dots, x_n) - f(x_1, \dots, x_i - h, \dots, x_n)}{2h},$$

for a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ , where  $h$  is a prescribed resolution parameter.

In total, five evaluations of  $\tilde{s}_N^M$  are required to approximate the derivatives in  $H_{R_N^M}$  and  $\nabla R_N^M$  at each Gauss-Newton iteration. The importance of rapid evaluation of the forward model  $\boldsymbol{\mu} \rightarrow \tilde{s}_N^M(\boldsymbol{\mu})$  is thus apparent, and indeed made possible by *i*) the reduced basis method with empirical interpolation and *ii*) the rapid output evaluation described in Section 6.2.5.

Theoretically speaking, the simple Gauss-Newton iteration scheme (6.76) is not very robust and will serve merely as a “proof of concept”. For example, we may experience iterates outside the feasible parameter domain, requiring a restart of the iterative process from (say) one of the edges, or oscillating iterates. Moreover, without any control of the step-length at each iteration, global convergence is not ensured. However, Gauss-Newton methods are often used to find the minimiser of least-squares problems [17], and in practice, we may expect the scheme to converge rapidly.

The reader is referred to [16] for the most recent development of an “uncertainty region” RB approach to time-dependent inverse problems. This method takes into account the numerical and measurement errors and may be used to construct a region  $\mathcal{P} \subset \mathcal{D}$  that is guaranteed to contain all parameters consistent with the measured data.

## 6.4 Numerical results

### 6.4.1 Spectral element truth approximation

For the parameter vectors  $\boldsymbol{\mu} = (1.5, 1.5)$ ,  $\boldsymbol{\mu} = (1.75, 1.5)$  and  $\boldsymbol{\mu} = (2, 2)$ , spectral element solutions to (6.17) are exhibited by Figure 6.5. Figures 6.5a, 6.5c and 6.5e show the solution with  $\epsilon_{\text{ano}} = 0.1$ , corresponding to an “insulating” anomaly, while Figures 6.5b, 6.5d and 6.5f show the solution with  $\epsilon_{\text{ano}} = 10$ , corresponding to a “superconductor” anomaly. In both cases, it is evident that the position of the anomaly affects the electric field lines and thus plausibly the electric flux across the electrode boundaries. As the background permittivity, we always choose  $\epsilon_{\text{bg}} = 1$ .

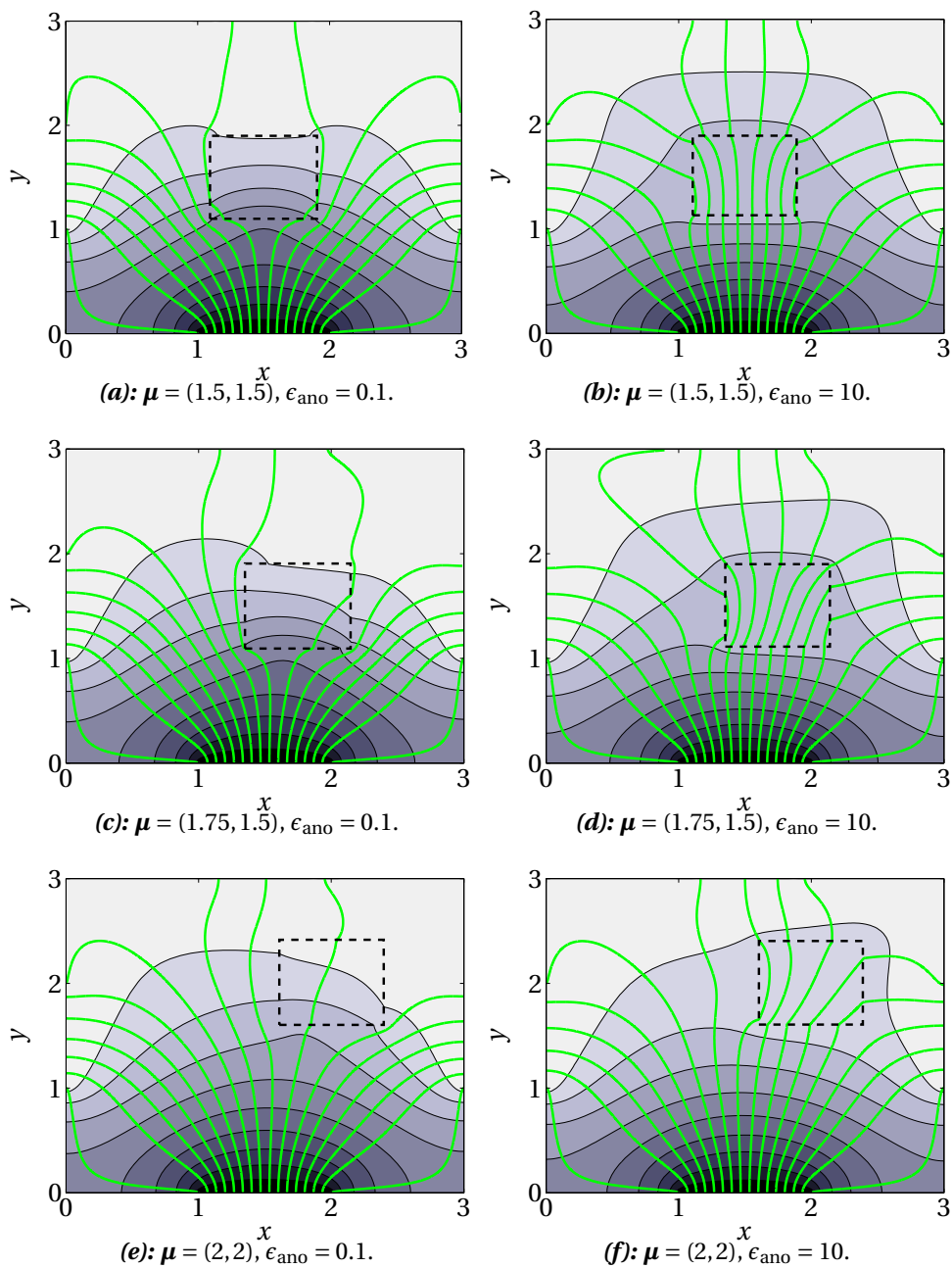
We now consider the numerical error in our spectral element approximation  $u_{\mathcal{N}} \approx u$ . In particular, we are interested in the order of convergence when we increase the polynomial degree  $P$  of the basis functions. As the exact solution is now not available for error comparison, we define the spectral element approximation  $u_{\mathcal{N}_e} \approx u$  as a surrogate for  $u$ . We assume that  $u_{\mathcal{N}_e}$  is computed with basis functions of a polynomial degree  $P_e$  sufficiently larger than  $P$  that  $|u_{\mathcal{N}} - u| \gg |u_{\mathcal{N}_e} - u|$ , and hence that

$$(6.78) \quad e_{\mathcal{N}}(\boldsymbol{\mu}) = u(\boldsymbol{\mu}) - u_{\mathcal{N}}(\boldsymbol{\mu}) \approx u_{\mathcal{N}_e}(\boldsymbol{\mu}) - u_{\mathcal{N}}(\boldsymbol{\mu}),$$

and

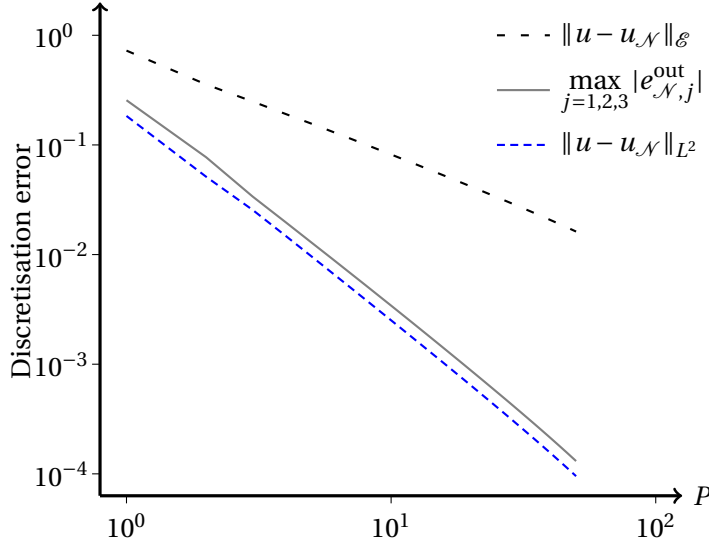
$$(6.79) \quad e_{\mathcal{N},j}^{\text{out}}(\boldsymbol{\mu}) = l_j^{\text{out}}(e_{\mathcal{N}}(\boldsymbol{\mu})) \approx l_j^{\text{out}}(u_{\mathcal{N}_e}(\boldsymbol{\mu}) - u_{\mathcal{N}}(\boldsymbol{\mu})),$$

## 6. A worked example: Electrostatics



**Figure 6.5:** Six truth approximations (spectral element solutions) for the problem (6.17) corresponding to different choices of the parameter vector  $\boldsymbol{\mu}$  and permittivity  $\epsilon_{\text{ano}}$ . The anomaly is either “superconductor-like”,  $\epsilon_{\text{ano}} = 10$ , or “insulator-like”,  $\epsilon_{\text{ano}} = 0.1$ . In all experiments,  $\epsilon_{\text{bg}} = 1$ . Electric field lines (green) and contour lines for the potential  $u_{\mathcal{N}_t}$  (background) are shown. The dashed square denotes the position of of the anomaly.





**Figure 6.6:** Discretisation errors as functions of the polynomial degree  $1 \leq P \leq 50$  for the spectral element solution to (6.17) with  $\boldsymbol{\mu} = \boldsymbol{\mu}_{\text{ref}} = (1.75, 1.5)$  and  $\epsilon_{\text{ano}} = 10$ . The plot shows algebraic convergence of order  $-2.08$  for the  $L^2$ -norm error and  $-1.03$  for the  $\mathcal{E}$ -norm error. For the maximum output error, the convergence is also of order  $-2.08$ .

for  $1 \leq j \leq 3$ . Below, we report results for  $e_{\mathcal{N}}$  and  $e_{\mathcal{N}}^{\text{out}}$ , even though the exact solution  $u$  is replaced with  $u_{\mathcal{N}_e}$  with  $P_e = 150$ .

In Figure 6.6, we show the energy-norm error  $\|e_{\mathcal{N}}(\boldsymbol{\mu})\|_{\mathcal{E}}$  and the  $L^2$ -norm error  $\|e_{\mathcal{N}}(\boldsymbol{\mu})\|_{L^2}$  in the field variable for the particular choice of parameter vector  $\boldsymbol{\mu} = \boldsymbol{\mu}_{\text{ref}} = (1.75, 1.5)$ , as well as the maximum error in the output, defined as

$$(6.80) \quad \max_{1 \leq j \leq 3} |e_{\mathcal{N},j}^{\text{out}}(\boldsymbol{\mu})|$$

for  $1 \leq P \leq 50$ . Clearly, the convergence is algebraic rather than exponential, suggesting that the exact solution is not analytic. We also note that the error in the output converges quadratically with the energy-norm error of the field variable, owing to the estimate (3.29), and that the decay in the  $L_2$ -error is of one order better than the decay in the energy error, which is as expected [3, 25].

In the model problem in Section 3.5.2, a non-convex corner in the physical domain was causing a singularity in the solution. The singularity was responsible for the limited algebraic convergence of the spectral element solution. In the present problem, the physical domain is a square and the domain is thus convex. However, strong singularities arise at the Dirichlet-Neumann interfaces at the edges of the electrodes, causing an even slower rate of convergence than what was achieved for the model problem.

Corner- and interface-singularities in solutions to elliptic partial differential equations are well-known phenomena [10, 25, 28]. One way to improve upon poor convergence rates due to the presence of singularities is refining the mesh close to the singularity. Another is the inclusion of special functions in the approximation space that mimic the behaviour of the singularities (see [7] for an example of successful application of such a technique to an electrostatics problem). A third is the *method of auxiliary mapping* [18], in which a small region containing the singularity is isolated and mapped onto a new domain through a mapping determined by *a priori* knowledge of the singular behaviour of the solution. The solution in the singular region is then approximated on the new domain using the existing basis functions. In this report, we do not pursue more efficient alternatives than the standard spectral element method, although such methods might unarguably prove beneficial.

### 6.4.2 Reduced basis approximation

We now present numerical results for our RB and RB-EI approximations. The same set of approximation spaces,  $X_N$ ,  $1 \leq N \leq N_{\max}$ , are used for both methods, and in all numerical experiments that follows,  $\epsilon_{\text{bg}} = 1$  and  $\epsilon_{\text{ano}} = 10$ .

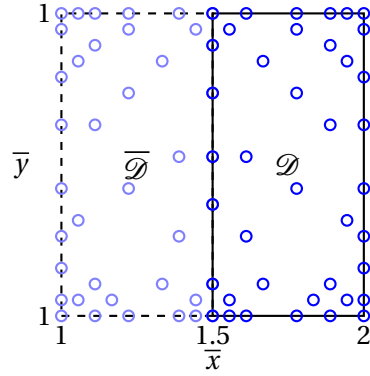
First, we consider the greedy sampling of parameter vectors. We then proceed with results from the RB approximation, including the *a posteriori* error estimators and the effect of different choices of the functions  $v_j^*$ . Finally, we apply the empirical interpolation method and examine the RB-EI approximation results.

#### Sampling procedure

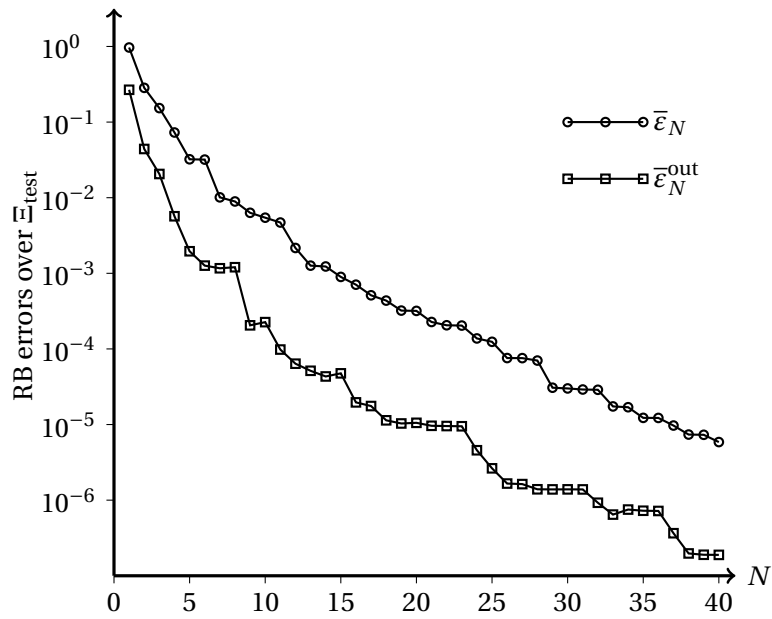
As our training sample  $\Xi_{\text{train}}$ , we choose a  $10 \times 20$  linear sample over  $\mathcal{D}$ . Relying on the output error bound  $\Delta_N^{\text{out,max}}$  defined in (6.72), the greedy algorithm (Algorithm 4.1) is then employed in the selection of  $S_N$ ,  $1 \leq N \leq N_{\max} = 40$ . Figure 6.7 shows the selected parameter vectors in  $\mathcal{D}$ , along with the points in the “mirrored” space  $\overline{\mathcal{D}}$ . We note that many of the vectors are chosen as extreme values, close to the edges and corners of the parameter space.

#### RB approximation results

Due to the greedy parameter selection, our reduced basis approximation will be accurate over the training sample. To test the approximation properties of

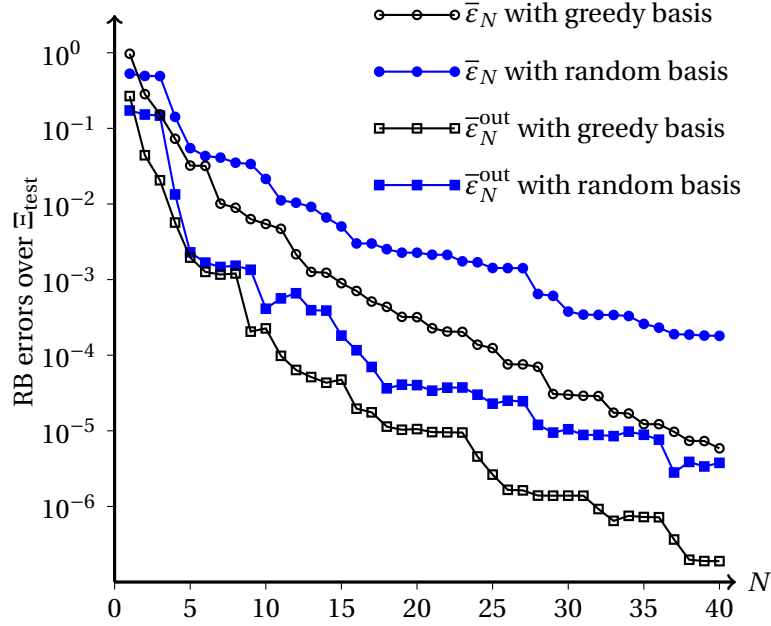


**Figure 6.7:** Choices of parameter vectors belonging to  $\mathcal{D}$  made by the greedy algorithm during construction of  $X_N$ . The “mirrored” parameter space  $\bar{\mathcal{D}}$  is also shown.



**Figure 6.8:** Maximum output error  $\bar{\epsilon}_N^{\text{out}}$  (squares) and maximum energy-norm error  $\bar{\epsilon}_N$  (circles) over  $\Xi_{\text{test}}$  for  $1 \leq N \leq 40$ .

## 6. A worked example: Electrostatics



**Figure 6.9:** Maximum output error  $\bar{\epsilon}_N^{\text{out}}$  (squares) and energy-norm error  $\bar{\epsilon}_N$  (circles) over  $\Xi_{\text{test}}$  when  $X_N$  is constructed randomly (filled marks) and greedily for  $1 \leq N \leq 40$

$X_N$  outside  $\Xi_{\text{train}}$ , we also introduce a test sample  $\Xi_{\text{test}} \neq \Xi_{\text{train}}$ , consisting of 200 randomly chosen points in  $\mathcal{D}$ . We also define the error measures

$$(6.81) \quad \bar{\epsilon}_N \stackrel{\text{def}}{=} \max_{\boldsymbol{\mu} \in \Xi_{\text{test}}} \|e_N(\boldsymbol{\mu})\|_{\mathcal{E}},$$

and

$$(6.82) \quad \bar{\epsilon}_N^{\text{out}} \stackrel{\text{def}}{=} \max_{\substack{\boldsymbol{\mu} \in \Xi_{\text{test}} \\ 1 \leq j \leq 3}} |(s_N(\boldsymbol{\mu}))_j - (s_{\mathcal{N}_t}(\boldsymbol{\mu}))_j|.$$

In Figure 6.8, we then show  $\bar{\epsilon}_N$  and  $\bar{\epsilon}_N^{\text{out}}$  for  $1 \leq N \leq N_{\text{max}}$ . Here, the outputs are computed through the bilinear form with the  $v_j^*$  chosen as the solutions of (6.44) with  $\boldsymbol{\mu}$  replaced by  $\boldsymbol{\mu}_{\text{ref}}$ . We note that the errors decay very rapidly compared to our achievements with the spectral element method and that the energy-norm error decays monotonically (as it must, by definition) whereas the output-error does not.

Instead of constructing  $X_N$  greedily, we could simply have chosen the parameters for which to compute snapshots randomly. In Figure 6.9, output and energy-norm errors over  $\Xi_{\text{test}}$  are compared for the greedy and random sampling strategies. While the greedy basis undoubtedly yields the best results,

$N$	$\eta_N^{\max}$	$\eta_N^{\text{med}}$	$\eta_N^{\text{mean}}$	$\eta_N^{\text{out,max}}$	$\eta_N^{\text{out,mean}}$	$\eta_N^{\text{out,med}}$
5	3.176	1.806	1.848	$2.398 \cdot 10^4$	$3.960 \cdot 10^2$	$1.212 \cdot 10^2$
10	2.894	1.795	1.825	$1.374 \cdot 10^4$	$3.201 \cdot 10^2$	$1.142 \cdot 10^2$
15	2.972	1.797	1.827	$1.885 \cdot 10^4$	$4.329 \cdot 10^2$	$1.337 \cdot 10^2$
20	2.804	1.791	1.810	$7.169 \cdot 10^4$	$7.653 \cdot 10^2$	$1.955 \cdot 10^2$
25	2.919	1.779	1.807	$9.354 \cdot 10^4$	$1.236 \cdot 10^3$	$3.047 \cdot 10^2$
30	2.631	1.762	1.789	$1.687 \cdot 10^5$	$1.020 \cdot 10^3$	$2.370 \cdot 10^2$
35	2.658	1.760	1.787	$1.406 \cdot 10^5$	$1.467 \cdot 10^3$	$3.363 \cdot 10^2$
40	2.765	1.771	1.802	$7.354 \cdot 10^4$	$9.966 \cdot 10^2$	$2.847 \cdot 10^2$

**Table 6.1:** Maximum, median and mean of the effectivities  $\eta_N(\boldsymbol{\mu})$  and  $\eta : N, i^{\text{out}}(\boldsymbol{\mu})$  for  $\boldsymbol{\mu} \in \Xi_{\text{test}}$  for a few values of  $N$

even the random basis performs quite well. This suggests that  $u(\boldsymbol{\mu})$  is very smooth in the parameter, as the approximation properties of  $X_N$  seem to be rather insensitive to the particular choice of snapshots.

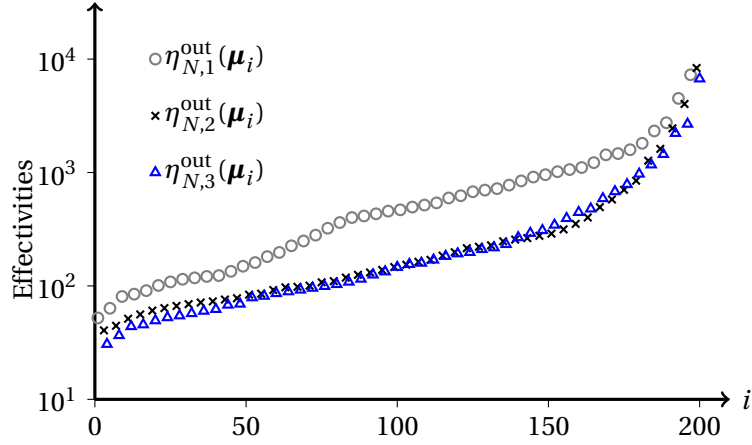
Although the *a posteriori* estimators  $\Delta_N$  and  $\Delta_{N,j}^{\text{out}}$  (see (6.69) and (6.71), respectively) are rigorous upper bounds, it is also important – for parameter selection as well as for certification of the RB solution – that they are quite sharp, meaning that the *effectivities*

$$(6.83) \quad \eta_N(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \frac{\Delta_N(\boldsymbol{\mu})}{\|e_N(\boldsymbol{\mu})\|_{\mathcal{E}}},$$

$$(6.84) \quad \eta_{N,j}^{\text{out}} \stackrel{\text{def}}{=} \frac{\Delta_{N,j}^{\text{out}}(\boldsymbol{\mu})}{\left| (\underline{s}_N(\boldsymbol{\mu}))_j - (\underline{s}_{\mathcal{N}_t}(\boldsymbol{\mu}))_j \right|},$$

where the latter is defined for  $1 \leq j \leq 3$ , are small. In Table 6.1 we present  $\eta_N^{\max}$ ,  $\eta_N^{\text{mean}}$  and  $\eta_N^{\text{med}}$ , denoting the maximum, mean and median of  $\eta_N(\boldsymbol{\mu})$  over  $\Xi_{\text{test}}$ , respectively, and  $\eta_N^{\text{out,max}}$ ,  $\eta_N^{\text{out,mean}}$  and  $\eta_N^{\text{out,med}}$ , denoting the maximum, mean and median of  $\eta_{N,j}^{\text{out}}(\boldsymbol{\mu})$ ,  $1 \leq j \leq 3$  over  $\Xi_{\text{test}}$ , respectively, for a few values of  $N$ . The effectivities corresponding to the energy-norm estimator are small, and so we can conclude that our coercivity lower bound  $\alpha_{\text{LB}}$  developed in Section 6.2.7 is quite sharp. Unfortunately, the effectivities for the output are quite large, (but not dramatically increasing with  $N$ ). However, we may hope that the  $\eta_N^{\text{out,max}}$  are rare deviations from the general behavior of  $\eta_N^{\text{out}}(\boldsymbol{\mu})$  over  $\mathcal{D}$ , and that we may for most  $\boldsymbol{\mu} \in \mathcal{D}$  expect  $\eta_N^{\text{out}}(\boldsymbol{\mu})$  to be more in compliance with the effectivity mean or median. In Figure 6.10, the  $\eta_{N,i}^{\text{out}}(\boldsymbol{\mu})$  are reported for  $\boldsymbol{\mu}$  across  $\Xi_{\text{test}}$  for the particular case  $N = 20$ . The data is sorted according to increasing effectivity. Roughly speaking, we observe that  $10^2 < \eta_{N,1}^{\text{out}} < 10^3$ , whereas  $\eta_{N,2}^{\text{out}}$  and

## 6. A worked example: Electrostatics



**Figure 6.10:** The output effectivities  $\eta_{N,1}(\boldsymbol{\mu})$ ,  $\eta_{N,2}(\boldsymbol{\mu})$  and  $\eta_{N,3}(\boldsymbol{\mu})$  for  $\boldsymbol{\mu} \in \Xi_{\text{test}}$  for the particular case  $N = 20$ . The effectivities are sorted in increasing order, and only one third is actually shown.

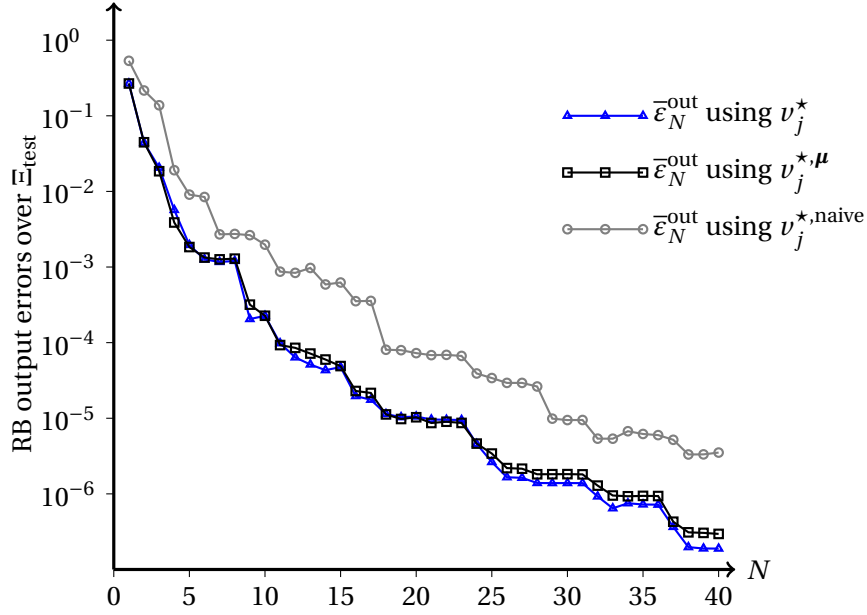
$\eta_{N,3}^{\text{out}}$  are slightly smaller. It seems plausible that the qualitative difference in the effectivities relates to the halving of the parameter space. At any rate, the output estimators must be regarded as quite conservative.

We now consider alternative choices of the output evaluation functions  $v_j^*$  corresponding to the alternatives described in Section 6.2.5:

- i)*  $v_j^*$  as the spectral element solution to (6.44) using approximation spaces of cubic polynomials. We denote this choice as  $v_j^{*,\boldsymbol{\mu}}$  due to the dependence upon the parameter vector. Note, however, that the computation of the  $v_j^{*,\boldsymbol{\mu}}$  are very fast due to the low polynomial order.
- ii)*  $v_j^*$  as the solution to (6.44) with  $\boldsymbol{\mu}$  replaced by  $\boldsymbol{\mu}_{\text{ref}}$  – thus the  $v_j^*$  can be computed as part of the preprocessing stage.
- iii)*  $v_j^*$  equal to unity on  $\Gamma_j$  and equal to zero in all other nodes. We denote this choice as  $v_j^{*,\text{naive}}$ .

The maximum output errors over  $\Xi_{\text{test}}$  for the three different choices of  $v_j^*$  are exhibited by Figure 6.11. We observe that alternative *i)* and *ii)* yields practically indistinguishable results, whereas alternative *iii)* results in output errors that are about one order of magnitude larger than the others.

Finally, we should emphasise that the smallest RB errors we here achieve are of little practical value, as the errors in the truth approximations are of order  $\sim 10^{-3}$  and thus much *larger* than the errors  $|u_{\mathcal{N}_t} - u|$ . Nevertheless, as a theo-



**Figure 6.11:** Maximum output error  $\bar{\epsilon}_N^{\text{out}}$  over  $\Xi_{\max}$  for different choices of  $v_j^*$  for output evaluation.

retical demonstration of the convergence properties of the method, the results are indeed very pleasing.

### RB-EI approximation results

We now apply the empirical interpolation method described in Chapter 5 to affinely approximate the bilinear form. Unfortunately, the gain of rapid computation comes at the price of an additional numerical error. Therefore, we now compare the error in the RB-EI solution and output to the error in the RB solution and output, respectively.

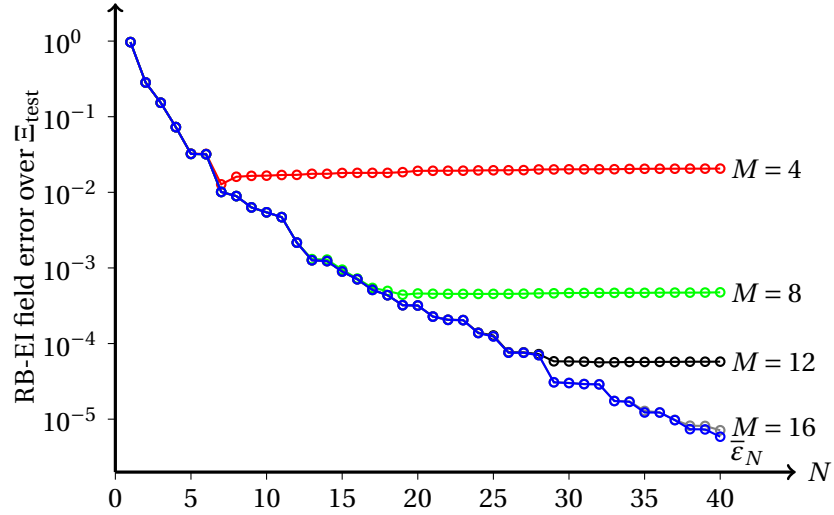
As RB-EI error measures, define

$$(6.85) \quad \bar{e}_N^M \stackrel{\text{def}}{=} \max_{\boldsymbol{\mu} \in \Xi_{\text{test}}} \|e_N^M(\boldsymbol{\mu})\|_{\mathcal{E}},$$

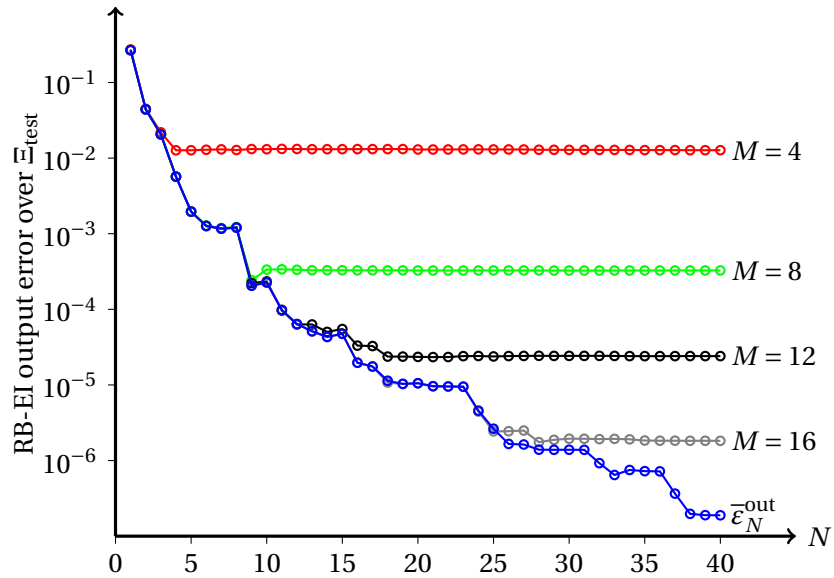
and

$$(6.86) \quad \bar{e}_N^{M,\text{out}} \stackrel{\text{def}}{=} \max_{\substack{\boldsymbol{\mu} \in \Xi_{\text{test}} \\ 1 \leq j \leq 3}} \left| (\tilde{s}_N^M(\boldsymbol{\mu}))_j - (s_{\mathcal{N}_t}(\boldsymbol{\mu}))_j \right|,$$

## 6. A worked example: Electrostatics



**Figure 6.12:** Maximum energy-norm errors  $\bar{e}_N^M$  over  $\Xi_{\text{test}}$ .



**Figure 6.13:** Maximum output errors  $\bar{e}_N^{M,\text{out}}$  over  $\Xi_{\text{test}}$ . The output is evaluated according to the procedure in Section 6.2.5.



where  $\underline{\tilde{s}}_N^M$  corresponds to the first row of the RB-EI output matrix, defined in (6.38).

As input to the empirical interpolation algorithm (Algorithm 5.1), we choose  $\Xi_{\text{train}}$  as the training sample (the same sample we use for construction of the RB space),  $M_{\text{max}} = 40$  and  $P_{\text{EI}} = P_t = 30$ .

In Figure 6.12 and 6.13, we show  $\bar{e}_N^M$  and  $\bar{e}_N^{M,\text{out}}$ , respectively, for  $1 \leq N \leq 40$  and a few values of  $M$ . We observe that the errors decay more or less identically to  $\bar{e}_N$  and  $\bar{e}_N^{\text{out}}$ , i.e. *without* empirical interpolation, until a certain point at which the interpolation error dominates completely. In fact, for  $M = 8$  and  $N \gtrsim 10$ , the error in the RB-EI output is of the same order as the error in the truth output.

### 6.4.3 Parameter estimation

#### Results with exact capacitance measurements

We now employ the inversion scheme described in Section 6.3. As a surrogate for the “exact” capacitance measurements, we use  $\tilde{s}_N^M(\boldsymbol{\mu})$  with  $N = 40$  and  $M = 40$ . We then “measure”  $s_{\text{obs}} \stackrel{\text{def}}{=} \tilde{s}_{40}^{40}(\boldsymbol{\mu}_{\text{exact}})$ , where  $\boldsymbol{\mu}_{\text{exact}} = (1.75, 1.5)$  is the exact position of the anomaly.

We then employ the procedure described in Section 6.3 with  $\tilde{s}_N^M(\boldsymbol{\mu})$  as the RB-EI output for the particular cases  $N \in \{4, 8, 12, 16\}$  and  $M = 16$ . We set  $h = 10^{-4}$  and stop the iterative process when the norm of the (Euclidian) distance between two successive iterates is less than  $10^{-5}$ .

The scheme converges in typically 5-6 iterations to  $\boldsymbol{\mu}_{N,\text{rec}}^M$ . As a measure of accuracy, we exhibit in Table 6.2 (middle column) the reconstruction errors

$$(6.87) \quad e_{N,\text{rec}}^M \stackrel{\text{def}}{=} \|\boldsymbol{\mu}_{N,\text{rec}}^M - \boldsymbol{\mu}_{\text{exact}}\|_2,$$

where  $\|\cdot\|_2$  denotes the standard Euclidian norm. We observe that even for small  $N$ , the reconstructed and exact anomaly positions are very close.

#### Results with noisy capacitance measurements

We now repeat the experiment above, but this time the capacitance measurements suffer from 1 percent (relative) white noise. Now denoting the reconstructed anomaly position by  $\boldsymbol{\mu}_{N,\text{rec},1\%}^M$ , the resulting reconstruction errors

$$(6.88) \quad e_{N,\text{rec},1\%}^M \stackrel{\text{def}}{=} \|\boldsymbol{\mu}_{N,\text{rec},1\%}^M - \boldsymbol{\mu}_{\text{exact}}\|_2,$$

## 6. A worked example: Electrostatics

---

$N$	$e_{N,\text{rec}}^M$	$e_{N,\text{rec},1\%}^M$
4	$1.56 \cdot 10^{-2}$	$2.68 \cdot 10^{-2}$
8	$4.73 \cdot 10^{-4}$	$1.39 \cdot 10^{-2}$
12	$1.15 \cdot 10^{-5}$	$1.35 \cdot 10^{-2}$
16	$2.04 \cdot 10^{-6}$	$1.35 \cdot 10^{-2}$

**Table 6.2:** Distances between the exact anomaly position  $\boldsymbol{\mu}_{\text{exact}} = (1.75, 1.5)$  and the reconstructed anomaly positions  $\boldsymbol{\mu}_{N,\text{rec}}^M$  (middle column) and  $\boldsymbol{\mu}_{N,\text{rec},1\%}^M$  (right column) for the particular case  $M = 16$ .

for  $N \in \{4, 8, 12, 16\}$  and  $M = 16$ , are given in the right column of Table 6.2.

We observe that there is no effect of increasing  $N$  beyond the point at which the inversion error is less than about  $10^{-2}$ , which seems reasonable due to the now noisy observations. In a real-life parameter-estimation problem though, a reconstruction error even of this order may be sufficiently accurate.

# Chapter 7

## Conclusions

The main results in this report relate to the evaluation of flux-type output functionals whose argument is the solution of a parametrised partial differential equation. In Section 2.5, we discussed a Neumann-Dirichlet equivalence allowing for convenient evaluation of the output through the bilinear form corresponding to the PDE at hand.

In Section 6.2.5, we showed how the output may be evaluated very efficiently through a slightly modified reduced basis stiffness matrix when either the problem at hand is affine or the empirical interpolation method is invoked. Moreover, the numerical experiments from Section 6.4.2 show that the output evaluation is very accurate as long as the number of (empirical) interpolation nodes,  $M$ , is not chosen too small.

The evaluation of the output through the bilinear form requires a particular function  $\nu^* \in V^*$  to be defined. If the argument of the output functional is a standard finite- or spectral element solution, every choice of  $\nu^*$  is equivalent, whereas this is not the case when the argument is a reduced basis solution. This has been theoretically argued for and supported by numerical results in Sections 2.5 and 6.4.2, respectively.

For the elaborate reduced basis example problem from Chapter 6, it remains to develop *a posteriori* error estimators that *i)* also incorporate the error resulting from the employment of empirical interpolation, and *ii)* allows for an offline-online computational approach, thus yielding very rapid online error bounds for the RB-EI output. By following a similar procedure as in [8, Chapter 6], this seems to be within reach both theoretically and practically. However, our investigation of the standard error estimators for the flux-type output of interest revealed they are very conservative, and so an effort should be put into the de-

## 7. Conclusions

---

velopment of sharper error bounds. In fact, by considering an adjoint problem (Section 4.4.2) for each output functional, not only may we accelerate the convergence of the RB outputs, but sharper *a posteriori* bounds for the RB outputs may be constructed based on the estimators for the primal and adjoint problems [20, 23].

We have also touched upon several topics which are in need of further investigation. Suggestions for future work include:

- the implementation of a more efficient method to deal with singularities that arise in our spectral element approximations, for instance by the method of auxiliary mapping [18],
- the treatment of inhomogeneous Dirichlet boundary conditions in the reduced basis context.

# Bibliography

- [1] Maxime Barrault, Yvon Maday, Ngoc Cuong Nguyen, and Anthony T. Patera. An ‘empirical interpolation’ method: application to efficient reduced-basis discretization of partial differential equations. *C. R. Math. Acad. Sci. Paris*, 339(9):667–672, 2004.
- [2] Christine Bernardi and Yvon Maday. Polynomial approximation of some singular functions. *Appl. Anal.*, 42(1):1–32, 1991.
- [3] Christine Bernardi and Yvon Maday. Spectral methods. In *Handbook of numerical analysis, Vol. V*, Handb. Numer. Anal., V, pages 209–485. North-Holland, Amsterdam, 1997.
- [4] C. Canuto and A. Quarteroni. Approximation results for orthogonal polynomials in Sobolev spaces. *Math. Comp.*, 38(157):67–86, 1982.
- [5] Gene H. Golub and Charles F. Van Loan. *Matrix computations*. Johns Hopkins Studies in the Mathematical Sciences. Johns Hopkins University Press, Baltimore, MD, third edition, 1996.
- [6] William J. Gordon and Charles A. Hall. Construction of curvilinear coordinate systems and applications to mesh generation. *Internat. J. Numer. Methods Engrg.*, 7:461–477, 1973.
- [7] D. Greenfield and M. Monastyrski. Three-dimensional electrostatic field calculation with effective algorithm of surface charge singularities treatment based on the Fichera’s theorem. *Nuclear Instruments and Methods in Physics Research, Section A*, 519:82–89, 2004.
- [8] M. A. Grepl. *Reduced-Basis Approximation and A Posteriori Error Estimation for Parabolic Partial Differential Equations*. PhD thesis, Massachusetts Institute of Technology, 2005.

## Bibliography

---

- [9] Martin A. Grepl, Yvon Maday, Ngoc C. Nguyen, and Anthony T. Patera. Efficient reduced-basis treatment of nonaffine and nonlinear partial differential equations. *M2AN Math. Model. Numer. Anal.*, 41(3):575–605, 2007.
- [10] P. Grisvard. *Singularities in boundary value problems*, volume 22 of *Recherches en Mathématiques Appliquées [Research in Applied Mathematics]*. Masson, Paris, 1992.
- [11] Erwin Kreyszig. *Introductory functional analysis with applications*. Wiley Classics Library. John Wiley & Sons Inc., New York, 1989.
- [12] John M. Lee. *Introduction to smooth manifolds*, volume 218 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, 2003.
- [13] Alf Emil Løvgren, Yvon Maday, and Einar M. Rønquist. A reduced basis element method for the steady Stokes problem. *M2AN Math. Model. Numer. Anal.*, 40(3):529–552, 2006.
- [14] Robert C. McOwen. *Partial Differential Equations: Methods and Applications*. Prentice Hall, New Jersey, 2003.
- [15] N. C. Nguyen. *Reduced-Basis Approximations and A Posteriori Error Bounds for Nonaffine and Nonlinear Partial Differential Equations: Application to Inverse Analysis*. PhD thesis, Singapore-MIT Alliance, National University of Singapore, 2005.
- [16] N. C. Nguyen, M. A. Grepl, A. T. Patera, and G. R. Liu. An "uncertainty region" reduced basis approach to parameter estimation for linear parabolic partial differential equations. Submitted to *Inverse Problems*. <http://augustine.mit.edu/>.
- [17] Jorge Nocedal and Stephen J. Wright. *Numerical optimization*. Springer Series in Operations Research. Springer-Verlag, New York, 1999.
- [18] Hae Soo Oh and Ivo Babuška. The  $p$ -version of the finite element method for the elliptic boundary value problems with interfaces. *Comput. Methods Appl. Mech. Engrg.*, 97(2):211–231, 1992.
- [19] A. T. Patera and G. Rozza. Reduced basis approximation and a posteriori error estimation for parametrized partial differential equations, version 1.0. <http://augustine.mit.edu>. Copyright MIT 2006–2007, to appear in (tentative rubric) *MIT Pappalardo Graduate Monographs in Mechanical Engineering*.

- 
- [20] A. T. Patera and E. M. Rønquist. A general output bound result: application to discretization and iteration error estimation and control. *Math. Models Methods Appl. Sci.*, 11(4):685–712, 2001.
- [21] Anthony T. Patera. A spectral element method for fluid dynamics: Laminar flow in a channel expansion. *Journal of computational physics*, 54:468–488, 1984.
- [22] Anthony T. Patera and Einar M. Rønquist. Reduced basis approximation and a posteriori error estimation for a Boltzmann model. *Computer Methods in Applied Mechanics and Engineering*, 196(29-30):2925–2942, 2007.
- [23] C. Prud’homme, D. V. Rovas, K. Veroy, L. Machiels, Y. Maday, A. T. Patera, and G. Turinici. Reliable real-time solution of parametrized partial differential equations: Reduced-basis output bound methods. *Journal of Fluids Engineering*, 124(1):70–80, 2002.
- [24] Alfio Quarteroni and Gianluigi Rozza. Numerical solution of parametrized Navier-Stokes equations by reduced basis methods. *Numer. Methods Partial Differential Equations*, 23(4):923–948, 2007.
- [25] Alfio Quarteroni and Alberto Valli. *Numerical approximation of partial differential equations*, volume 23 of *Springer Series in Computational Mathematics*. Springer-Verlag, Berlin, 1994.
- [26] Einar M. Rønquist. Numerical solution of partial differential equations. Lecture notes (course MA8502), NTNU, 2007.
- [27] Julius A. Stratton. *Electromagnetic Theory*. McGraw-Hill, New York, 1941.
- [28] Barna Szabó and Ivo Babuška. *Finite element analysis*. A Wiley-Interscience Publication. John Wiley & Sons Inc., New York, 1991.

