



# Reduced-basis solutions of affinely-parametrized linear partial differential equations

Alexandre Janon

## ► To cite this version:

Alexandre Janon. Reduced-basis solutions of affinely-parametrized linear partial differential equations. [Research Report] LJK. 2010. hal-00725344

**HAL Id: hal-00725344**

**<https://hal.science/hal-00725344>**

Submitted on 24 Aug 2012

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Reduced-basis solutions of affinely-parametrized linear partial differential equations

Alexandre Janon<sup>1</sup>

June 2010

---

## Contents

<b>1</b>	<b>Notations and assumptions</b>	<b>3</b>
<b>2</b>	<b>Reduced basis approximation</b>	<b>4</b>
2.1	Opportunity for a RB approximation . . . . .	4
2.2	Offline and online procedures . . . . .	4
<b>3</b>	<b>Error bound</b>	<b>6</b>
3.1	Derivation of the error bound . . . . .	6
3.2	Dual norm of the residual $\ r_\mu\ _{X'}$ . . . . .	8
3.3	Stability constant $\alpha(\mu)$ . . . . .	10
3.3.1	An easy case . . . . .	10
3.3.2	The general case: successive constraints method (SCM) . . . . .	11
<b>4</b>	<b>Choice of a reduced basis</b>	<b>15</b>
4.1	Greedy procedure . . . . .	15
4.2	Proper Orthogonal Decomposition (POD)-driven procedure . . . . .	16
<b>5</b>	<b>Case of a scalar output</b>	<b>21</b>
5.1	Primal approach . . . . .	21
5.2	Primal/dual approach . . . . .	22
5.3	Output-aware choice of basis: goal-oriented POD . . . . .	24

---

<sup>1</sup>Laboratoire Jean Kuntzmann, Université Grenoble 1/INRIA,  
<mailto:alexandre.janon@imag.fr>

## Introduction and overview

In some contexts (for example, in Monte-Carlo, or quasi-Monte-Carlo sensitivity analysis), it is necessary to compute the solution  $u_\mu$  of some partial differential equation (PDE) depending on an input parameter vector  $\mu$  (dependence can be obviously through the PDE coefficients, and also through the domain on which the equation is posed, or through initial/boundary data), for a large number (say, thousands) values of  $\mu$ . When no analytical solution of the PDE is known, one has to use computationally-intensive numerical methods – such as finite elements or finite differences – to approximate the solution of the PDE. This leads to complex computer codes, which can take several hours to run – for a *single* value of  $\mu$ . In this case, one easily sees that simply calling the numerical code for each required value of  $\mu$  – one-by-one – is not practicable at all ; and that we need to use a computationally cheaper code giving a *surrogate* solution  $\tilde{u}_\mu$  (also called *metamodel*) from a value of  $\mu$ . More specifically, we split our computation in two phases:

- an *offline* phase, where we make a reasonable number of calls to the expensive code in order to "learn" about the solutions of the PDE ;
- an *online* phase, where, given some value of  $\mu$ , we use the information collected during the offline phase to produce a good approximation of the solution of the PDE for this specific value of  $\mu$ , faster than if we had no information at all.

This strategy is efficient if the fixed additional cost of data collection during the offline phase (done once for every value of  $\mu$ ) is dominated by a strong reduction in the marginal cost (i.e., per value of  $\mu$ ) when one replaces the original, expensive numerical code by the cheaper "online" code.

In this report, we describe a procedure, named *reduced basis* approximation (RB), specifying the offline and online phases. The key to RB approximation is to look for (during the online phase) a surrogate  $\tilde{u}_\mu$  solution of the PDE as a linear combination of functions from a linearly-independent family (called the reduced basis) that has to be wisely chosen during the offline phase. For the resulting problem to be well-posed, one has to somewhat weaken the PDE, using a Galerkin-like method. An advantage of using such a procedure is that one can quantify the loss of information made when replacing the expensive numerical solution  $u_\mu$  by  $\tilde{u}_\mu$ , using an *a posteriori* fastly-computed *error bound*

$\varepsilon_\mu$  so that

$$\|u_\mu - \tilde{u}_\mu\| \leq \varepsilon_\mu$$

for some functional norm  $\|\cdot\|$ .

The report is organised as follows: in a first part, we present the necessary assumptions on the PDE for RB to be applicable to ; for the sake of simplicity, we restrict ourselves to the case of coercive, linear and "affinely parametrized" (we will explain what this means) PDEs, but it should be noted there are extensions to RB which somewhat relax some of these assumptions ; in a second part, we describe the core of the RB approximation, that are the offline and online phases of computation of the surrogate  $\tilde{u}_\mu$  ; in a third part, we give a procedure to compute the error bound  $\varepsilon_\mu$  ; in a fourth part, we show two different ways (*greedy* and *POD*) to choose the reduced basis family ; in a fifth and final part, we address the question of using RB to approximate *outputs* of the solution, that are (scalar) functionals of  $u_\mu$ .

## 1 Notations and assumptions

Our main reference for sections 1 to 3 is [5].

As said earlier, we denote by  $\mu \in \mathcal{P} \subset \mathbf{R}^p$  our parameter vector.

Our unknown is a function  $u_\mu$  — depending on  $\mu$ , besides the usual (e.g. time and space) variables — belonging to some function space  $X$ , satisfying a  $\mu$ -parametrized linear partial differential equation, which can be written under a variational (a.k.a. weak) form:

$$a_\mu(u_\mu, v) = f_\mu(v) \quad \forall v \in X \quad (1)$$

where  $a_\mu$  is a  $\mu$ -dependent (continuous) bilinear form on  $X$ , and  $f$  is a (continuous) linear form on  $X$ .

We further assume that our PDE is *affinely parametrized*, that is:

- $f_\mu \equiv f$  is  $\mu$ -independent;
- and  $a_\mu$  can be written as a sum of  $Q$  bilinear forms:

$$a_\mu(w, v) = \sum_{q=1}^Q \Theta_q(\mu) a_q(w, v)$$

where:

- $a_q$  are  $\mu$ -independent bilinear forms on  $X$  ;
- $\Theta_q$  are arbitrary functions on  $\mathcal{P}$ , capturing all dependence in  $\mu$  ; since these functions are to be evaluated many times during the online phase, we shall require these functions to be fast-computed.

## 2 Reduced basis approximation

### 2.1 Opportunity for a RB approximation

The function space  $X$  is typically a Sobolev space; such spaces, due to their infinite-dimensional nature, are not amenable to numerical computations. In order to discretize equation (1) and get an approximation of the actual mathematical solution, one can set  $X$  to some finite-element space (e.g.,  $\mathbb{P}^1$ ) and use (1) to get a linear system of equations satisfied by the coefficients of  $u_\mu$  in a basis of  $X$ . In the sequel,  $X$  will always stand for such a space. This means that the reduced basis solution we are to compute is an approximation *to the discrete numerical solution*, and thus we rely on (and assume) the fact that the latter is a good approximation to the analytical solution. A good treatment of finite element discretizations can be found in [7].

The computational cost required to compute this approximation is an increasing function of the dimension of the chosen discrete function space, and unfortunately, it has to be very large for the approximation to be accurate. One explanation for this fact is that those spaces are "too general", that is, they contain too many different functions and so we need to deal with a large number of degrees of freedom for our unknown.

On the other hand, we expect that two solutions of (1), for two different values of  $\mu$ , will exhibit common "features". We could "encode" those features into a space  $\tilde{X}$ , spanned by some ( $\mu$ -independent) functions  $\zeta_1, \zeta_2, \dots, \zeta_N \in X$ , and replace  $X$  by  $\tilde{X}$  in (1). The family  $\{\zeta_1, \dots, \zeta_N\}$  is called the *reduced basis*. For a sufficiently regular-in- $\mu$  problem, we expect that taking  $N \ll \dim X$  will still give – for a well-choosed reduced basis – an accurate enough solution, yet providing substantial computational economy.

### 2.2 Offline and online procedures

We suppose, for now, that we are given a reduced basis  $\{\zeta_1, \dots, \zeta_N\}$  of linearly independent functions in  $X$  (discussion about choice of this basis can be found

in section 4). We set  $\tilde{X} = \text{Span}\{\zeta_1, \dots, \zeta_N\}$ .

Given  $\mu \in \mathcal{P}$ , we look for  $\tilde{u}_\mu$  as a function belonging to  $\tilde{X}$ , thus :

$$\tilde{u}_\mu = \sum_{n=1}^N c_n(\mu) \zeta_n \quad (2)$$

where the unknowns are  $c_n(\mu) \in \mathbf{R}$ ,  $n = 1, \dots, N$ .

These unknowns are to be found by imposing that  $\tilde{u}_\mu$  satisfy a relaxed version of (1); instead of verifying (1) for every  $v \in X$ ,  $\tilde{u}_\mu$  satisfies (1) for every  $v \in \tilde{X}$ :

$$a_\mu(\tilde{u}_\mu, v) = f(v) \quad \forall v \in \tilde{X} \quad (3)$$

Since  $\tilde{X} \subset X$ , (3) is weaker than (1). This standard technique is known as *Galerkin projection*.

Substituting (2) into (3), and setting  $v = \zeta_p$  ( $p = 1, \dots, N$ ) gives that the  $c_n(\mu)$ 's satisfy the following system of  $N$  linear equations:

$$\sum_{n=1}^N \left[ \sum_{q=1}^Q \Theta_q(\mu) a_q(\zeta_n, \zeta_p) \right] c_n(\mu) = f(\zeta_p) \quad \forall p = 1, \dots, N \quad (4)$$

This justifies the following offline/online procedure for the computation of  $\tilde{u}_\mu$ :

- in the offline phase:
  - we choose a reduced basis  $\{\zeta_1, \dots, \zeta_N\}$  (more details about this choice later)
  - we assemble, and store,  $Q$  matrices of size  $N \times N$ :

$$\mathbf{a}_q = \left( a_q(\zeta_n, \zeta_p) \right)_{p,n=1, \dots, N}$$

and the  $N$ -sized vector:

$$\mathbf{f} = \left( f(\zeta_p) \right)_{p=1 \dots N}$$

- in the online phase, being given a value of  $\mu$ :

- we assemble the  $N \times N$  matrix:

$$\mathbf{a}(\mu) = \sum_{q=1}^Q \Theta_q(\mu) \mathbf{a}_q$$

- our vector of unknowns  $\mathbf{c}(\mu) = (c_n(\mu))_{n=1\dots N}$  is given by :

$$\mathbf{c}(\mu) = (\mathbf{a}(\mu))^{-1} \mathbf{f}$$

where the matrix inversion solves the system (4),

- and so we get:

$$\tilde{u}_\mu = \sum_{n=1}^N c_n(\mu) \zeta_n$$

The complexity of the online phase is as follows: assembling of  $\mathbf{a}(\mu)$  takes  $O(QN^2)$  operations. Its inversion takes  $O(N^3)$  operations (one should note that  $\mathbf{a}(\mu)$  has no reason to be a sparse matrix – it is not, indeed). So the coefficients of  $\tilde{u}_\mu$  with respect to the reduced basis can be computed in less than  $O(QN^2 + N^3)$  operations.

It is important to note that this complexity is *independent* of  $\dim X$ . Thus online computation time is independent of the "quality" of the underlying reference numerical approximation, which can be as good as wanted without impacting marginal cost. We notice that we neglected the cost of evaluations of  $\Theta_q(\mu)$ , hence the hypothesis of fast-evaluation of  $\Theta_q$  we have made. We also notice that online complexity (in computational time and in storage requirements) is proportional to  $Q$ , the number of terms in the affine decomposition of  $a_\mu$ ; for our procedure to be efficient,  $Q$  should therefore be kept orders of magnitude smaller than  $\dim X$ .

### 3 Error bound

#### 3.1 Derivation of the error bound

In this section, we design an error bound  $\varepsilon_\mu$  so that:

$$\|u_\mu - \tilde{u}_\mu\| \leq \varepsilon_\mu$$

for all  $\mu \in \mathcal{P}$ .

We should note once again that, thanks to our assumption about the quality of the underlying approximation, we measure only the error between the reference, expensive numerical solution  $u$  and our surrogate  $\tilde{u}$ ; error between  $u$  and the analytical solution is neglected.

We endow  $X$  with an Hilbert structure: an inner product  $\langle \cdot, \cdot \rangle$  on  $X \times X$

yielding a norm  $\|\cdot\|$  on  $X$ . We need a crucial assumption about the bilinear form  $a_\mu$  describing our PDE: the  $\alpha(\mu)$  defined by:

$$\alpha(\mu) = \inf_{w \in X} \sup_{v \in X} \frac{a_\mu(w, v)}{\|w\| \|v\|}$$

has to be strictly positive:  $\alpha(\mu) > 0 \quad \forall \mu \in \mathcal{P}$ .

The quantity  $\alpha(\mu)$  is known as Babuška's "inf-sup" stability constant. The inf-sup stability hypothesis  $\alpha(\mu) > 0$  is a sufficient condition for well-posedness of the underlying variational problem (1) (Babuška-Brezzi condition).

In the sequel, we suppose that  $a_\mu$  is *symmetric*. The inf-sup stability hypothesis is then equivalent to the coercivity of  $a_\mu$  (and the well-posedness of the variational problem comes from Lax-Milgram), and we have:

$$\alpha(\mu) = \inf_{w \in X} \frac{a_\mu(w, w)}{\|w\|^2}$$

Let's move on to the error bound. First we remember that:

$$a_\mu(u_\mu, v) = f_\mu(v) \quad \forall v \in X$$

Subtracting  $a_\mu(\tilde{u}_\mu, v)$  from both sides of this equation, we get, thanks to bilinearity of  $a_\mu$ :

$$a_\mu(u_\mu - \tilde{u}_\mu, v) = f_\mu(v) - a_\mu(\tilde{u}_\mu, v) \quad \forall v \in X$$

So:

$$a_\mu(e_\mu, v) = r_\mu(v) \quad \forall v \in X \tag{5}$$

if we set:

$$e_\mu = u_\mu - \tilde{u}_\mu$$

$$r_\mu(v) = f(v) - a_\mu(\tilde{u}_\mu, v)$$

( $e_\mu$  is known as the *error*, and  $r_\mu$  as the *residual form*).

Introducing  $\alpha(\mu)$ , we have the following inequality:

$$a_\mu(e_\mu, e_\mu) \geq \alpha(\mu) \|e_\mu\|^2$$

We also have:

$$\|r_\mu(v)\| \leq \|r_\mu\|_{X'} \|v\|$$



where

$$\|r_\mu\|_{X'} = \sup_{v \in X, \|v\|=1} r_\mu(v)$$

is the dual norm of the residual.

Thanks to these inequalities, and (5) with  $v = e_\mu$ , we get:

$$\alpha(\mu)\|e_\mu\|^2 \leq \|r_\mu\|_{X'}\|e_\mu\|$$

and so, thanks to our stability assumption:

$$\|e_\mu\| = \|u_\mu - \tilde{u}_\mu\| \leq \frac{\|r_\mu\|_{X'}}{\alpha(\mu)} \quad (6)$$

So we have:

$$\varepsilon_\mu = \frac{\|r_\mu\|_{X'}}{\alpha(\mu)}$$

In order to practically exploit the error bound (6), it remains to explain how to efficiently work out  $\|r_\mu\|$  and  $\alpha(\mu)$ . This question is addressed in the two following sections.

### 3.2 Dual norm of the residual $\|r_\mu\|_{X'}$

Riesz representation theorem ensures that there exist  $y_\mu \in X$  so that

$$\langle y_\mu, v \rangle = r_\mu(v) \quad \forall v \in X \quad (7)$$

Moreover, computing  $y_\mu$  is sufficient because of the additional property:

$$\|r_\mu\|_{X'} = \|y_\mu\|$$

Equation (7) can be rewritten, thanks to the definition of  $r_\mu$  and the affine decomposition of  $a_\mu$ :

$$\langle y_\mu, v \rangle = f(v) - \sum_{q=1}^Q \Theta_q(\mu) a_q(\tilde{u}_\mu, v) \quad \forall v \in X$$

and so, making use of the reduced basis decomposition of  $\tilde{u}_\mu$  (equation (2)):

$$\langle y_\mu, v \rangle = f(v) - \sum_{q=1}^Q \sum_{n=1}^N c_n(\mu) a_q(\zeta_n, v) \quad \forall v \in X$$

This equation can be viewed as a linear variational problem whose unknown is  $y_\mu \in X$ . By Duhamel superposition principle,  $y_\mu$  can be written as:

$$y_\mu = \gamma - \sum_{q=1}^Q \sum_{n=1}^N \Theta_q(\mu) c_n(\mu) \Gamma_{qn} \quad (8)$$

where  $\gamma \in X$ , and  $\Gamma_{qn} \in X$  ( $q = 1, \dots, Q$ ,  $n = 1, \dots, N$ ) satisfy:

$$\langle \gamma, v \rangle = f(v) \quad \forall v \in X \quad (9)$$

$$\langle \Gamma_{qn}, v \rangle = a_q(\zeta_n, v) \quad \forall v \in X \quad (10)$$

Law of cosines applied to (8) gives:

$$\begin{aligned} \|y_\mu\|^2 &= \|\gamma\|^2 - 2 \sum_{q=1}^Q \sum_{n=1}^N \Theta_q(\mu) c_n(\mu) \langle \gamma, \Gamma_{qn} \rangle \\ &\quad + \sum_{q=1}^Q \sum_{q'=1}^Q \sum_{n=1}^N \sum_{n'=1}^N \Theta_q(\mu) c_n(\mu) \Theta_{q'}(\mu) c_{n'}(\mu) \langle \Gamma_{qn}, \Gamma_{q'n'} \rangle \end{aligned} \quad (11)$$

Equation (11) will be our key to an offline/online procedure for computing  $\|r_\mu\|_{X'}$ . Each of those phases should be done after the phases of computation of  $\tilde{u}(\mu)$  described in section 2.2. The procedures read:

- in the offline phase:
  - we solve the  $QN + 1$  linear variational problems (9) and (10) for  $\gamma$  and  $\Gamma_{qn}$  ( $q = 1, \dots, Q$ ,  $n = 1, \dots, N$ ) ;
  - we compute and store  $1 + QN + (QN)^2$  inner products:

$$\langle \gamma, \gamma \rangle, \langle \gamma, \Gamma_{qn} \rangle, \langle \Gamma_{qn}, \Gamma_{q'n'} \rangle$$

$$(q, q' = 1, \dots, Q, n, n' = 1, \dots, N)$$

- in the online phase, we compute  $\|r_\mu\|_{X'} = \|y_\mu\|$  using equation (11).

We note that the vectors  $\gamma$  and  $\Gamma_{qn}$  computed during offline phase are *not* saved for the online phase, we store only  $1 + QN + (QN)^2$  scalars. Thus in the online phase, the actual vector  $y_\mu$  is never formed. Doing so would make complexity of the online phase depend on  $\dim X$ ; this is something we would like to avoid.

### 3.3 Stability constant $\alpha(\mu)$

We now turn to the problem about estimating  $\alpha(\mu)$ .

Recall that

$$\alpha(\mu) = \inf_{v \in X, \|v\|=1} a_\mu(v, v)$$

$\alpha(\mu)$  is thus the optimal value of a minimization problem posed on  $X$ . The objective function in this problem is a quadratic form, we seek to minimize on the unit sphere. It is a well-known fact that  $\alpha(\mu)$  is the smallest eigenvalue of a (symmetric) matrix associated to the quadratic form  $a_\mu(\cdot, \cdot)$ . Indeed, if for all  $v$  and  $w$  in  $X$ :  $a_\mu(v, w) = v^T M_\mu w$  and if  $\langle v, w \rangle = v^T \Omega w$  for symmetric matrices  $M_\mu$  and  $\Omega$ , then  $\alpha(\mu)$  is the smallest  $\lambda \in \mathbf{R}$  so that equation :

$$M_\mu v = \lambda \Omega v$$

admits a nonzero solution  $v$ .

Due to the high-dimensionality of  $X$ , it would be too expensive to solve this (generalized) eigenproblem online, for each value of  $\mu$ . However, we expect that affine decomposition of  $a_\mu$  will provide us with an efficient offline/online procedure for estimating  $\alpha(\mu)$ . By estimating here we mean finding a lower bound for  $\alpha(\mu)$ , because the form of the error estimator (6) shows us we only need such a lower bound:

$$\tilde{\alpha}(\mu) \leq \alpha(\mu)$$

because we have therefore:

$$\|u - \tilde{u}\| \leq \frac{\|r\|_{X'}}{\tilde{\alpha}(\mu)}$$

#### 3.3.1 An easy case

We should first address the case where  $\Theta_q(\mu) \geq 0 \quad \forall q = 1, \dots, Q, \forall \mu \in \mathcal{P}$ .

If this condition is met, then for any  $\bar{\mu} \in \mathcal{P}, \mu \in \mathcal{P}$ :

$$\begin{aligned}
\sum_{q=1}^Q \Theta_q(\mu) a_q(u, v) &\geq \sum_{q=1}^Q \frac{\Theta_q(\mu)}{\Theta_q(\bar{\mu})} \Theta_q(\bar{\mu}) a_q(u, v) \\
&\geq \left\{ \inf_{q'=1, \dots, Q} \frac{\Theta_{q'}(\mu)}{\Theta_{q'}(\bar{\mu})} \right\} \sum_{q=1}^Q \Theta_q(\bar{\mu}) a_q(u, v) \\
&\geq \left\{ \inf_{q'=1, \dots, Q} \frac{\Theta_{q'}(\mu)}{\Theta_{q'}(\bar{\mu})} \right\} \alpha(\bar{\mu}) \|u\| \|v\|
\end{aligned}$$

This leads us to the following offline/online procedure for  $\tilde{\alpha}(\mu)$ :

- Offline:
  - choose a "nominal" value  $\bar{\mu} \in \mathcal{P}$  for which  $\Theta_q(\bar{\mu}) \neq 0 \ \forall q = 1, \dots, Q$ ;
  - compute and store  $\alpha(\bar{\mu})$  by solving an eigenproblem on  $X$ .
- Online, given a value of  $\mu$ , take

$$\tilde{\alpha}(\mu) = \left\{ \inf_{q=1, \dots, Q} \frac{\Theta_q(\mu)}{\Theta_q(\bar{\mu})} \right\} \alpha(\bar{\mu})$$

This procedure is simple and efficient, with a  $\dim X$ -independent online stage. However, if our positivity hypothesis on  $\Theta_q$  is violated, we have to resort to a more complex method, which we describe below.

### 3.3.2 The general case: successive constraints method (SCM)

The successive constraints method (SCM) is presented in [3].

$\alpha(\mu)$  can be rewritten

$$\alpha(\mu) = \inf_{w \in X, \|w\|=1} a_\mu(w, w)$$

We further note that:

$$\alpha(\mu) = \inf_{w \in X, \|w\|=1} \left\{ \sum_{q=1}^Q \Theta^q(\mu) a_q(w, w) \right\} = \inf_{y \in \mathcal{Y}} \left\{ \sum_{q=1}^Q \Theta^q(\mu) y_q \right\} \quad (12)$$

where:

$$\mathcal{Y} = \{(y_1, \dots, y_Q) \in \mathbf{R}^Q \text{ st. } \exists w \in X, \|w\| = 1 \text{ and } y_q = a_q(w, w) \forall q = 1, \dots, Q\}$$

To obtain a lower bound on  $\alpha(\mu)$ , we should replace  $\mathcal{Y}$  by a  $\tilde{\mathcal{Y}} \supseteq \mathcal{Y}$ ; in other words, we should relax the constraints on the optimization problem (12).

To define  $\tilde{\mathcal{Y}}$ , we begin by noticing that  $\mathcal{Y}$  is bounded: let:

$$\sigma_q^- = \inf_{w \in X, \|w\|=1} a_q(w, w), \quad \sigma_q^+ = \sup_{w \in X, \|w\|=1} a_q(w, w), \quad \mathcal{B} = \prod_{q=1}^Q [\sigma_q^-, \sigma_q^+].$$

Now let  $y = (y_1, \dots, y_Q) \in \mathcal{Y}$  and let  $w$  satisfying the properties in the definition of  $\mathcal{Y}$ . For every  $\mu' \in \mathcal{P}$ , we have:

$$\begin{aligned} \sum_q \Theta_q(\mu') y_q &= \sum_q \Theta_q(\mu') a_q(w, w) \\ &\geq \inf_{w' \in X, \|w'\|=1} \sum_q \Theta_q(\mu') a_q(w', w') \\ &= \alpha(\mu') \end{aligned}$$

So far we have:

$$\mathcal{Y} \subset \left\{ y \in \mathcal{B} \mid \sum_{q=1}^Q \Theta_q(\mu') y_q \geq \alpha(\mu') \quad \forall \mu' \in \mathcal{P} \right\}$$

The right-hand side of this inclusion would be a good candidate for  $\tilde{\mathcal{Y}}$ ; however it is defined by an infinite number of  $\mu$ -independent constraints (each for one  $\mu' \in \mathcal{P}$ ). We should relax one more time these constraints in order to get a finite number of (now,  $\mu$ -dependent) constraints. We take:

$$\tilde{\mathcal{Y}} = \left\{ y \in \mathcal{B} \mid \sum_{q=1}^Q \Theta_q(\mu') y_q \geq \alpha(\mu') \quad \forall \mu' \in \mathcal{S}_M(\mu, \mathcal{C}) \right\}$$

Where  $M \in \mathbf{N}$ ,  $\mathcal{C}$  is a finite subset of  $\mathcal{P}$  and  $\mathcal{S}_M(\mu, \mathcal{C})$  is the set of the  $M$  points in  $\mathcal{C}$  that are closest (with respect to some metric) to  $\mu$ .

So our SCM offline/online procedure is:

- offline:

- choose  $M$ ; choose  $\mathcal{C}$  (we postpone discussion about the choice of  $\mathcal{C}$  until the next subsection) ;
- for each  $\mu' \in \mathcal{C}$ , compute and store  $\alpha(\mu')$  by solving an eigenproblem on  $X$ ; compute and store  $\sigma_q^-$  and  $\sigma_q^+$  the same way, for each  $q = 1, \dots, Q$ .
- online: solve optimization problem:

$$\tilde{\alpha}(\mu) = \inf_{y=(y_1, \dots, y_Q) \in \tilde{\mathcal{Y}}} \left\{ \sum_{q=1}^Q \Theta^q(\mu) y_q \right\}$$

The online optimization problem is what is called a *linear programming problem*: the goal is to minimize a linear function under a finite number of linear inequalities constraints. In our case there are  $Q$  variables and  $Q+M$  constraints. There exist efficient algorithms, such as the simplex algorithm (see [6] for instance), which solve such optimization problems under (on average) polynomial complexity with respect to the number of variables and number of constraints, even if they can be exponential in the worst cases. The key here, once again, is that online complexity is still independent of  $\dim X$ .

A last remark we can do on the algorithm is about the trade-off in the choice of  $M$ : whatever  $M$  is, we always get a certified bound on  $\alpha(\mu)$ , but increasing  $M$  will improve sharpness of this bound, at the expense of an increase in online computation time.

#### Choice of $\mathcal{C}$

We will use a so-called *greedy* procedure for choosing the finite subset  $\mathcal{C} \subset \mathcal{P}$ . Roughly speaking, a decision algorithm is said to be *greedy* if he makes the best possible choice at each step. This strategy does not always yield a globally optimal decision but often gives a good solution for cases in direct global optimization would be unpractical.

For our greedy choice of  $\mathcal{C}$ , we will start with  $\mathcal{C} = \{\mu_1\}$ , with an arbitrarily chosen  $\mu_1 \in \mathcal{P}$ . Then, at each step of the algorithm, we are going to add a point in  $\mathcal{C}$  that gives the best possible expected improvement in the precision of our bound  $\tilde{\alpha}(\mu)$ .

Thus we get the following procedure for choice of  $\mathcal{C}$ :

1. choose  $M$  ;
2. initialize  $\mathcal{C} = \{\mu_1\}$  with some arbitrary  $\mu_1 \in \mathcal{P}$  ;
3. choose a finite-sized sample set  $\Xi \subset \mathcal{P}$  ;
4. repeat:
  - add

$$\mu^* = \underset{\mu \in \Xi}{\operatorname{argmax}} \alpha(\mu) - \tilde{\alpha}(\mu)$$

to  $\mathcal{C}$  (here  $\tilde{\alpha}(\mu)$  stands for the SCM lower bound computed using the "current"  $\mathcal{C}$ )

The repeat loop can be stopped either when  $\#\mathcal{C}$  has reached a maximal value, or when  $\max_{\mu \in \Xi} \alpha(\mu) - \tilde{\alpha}(\mu)$  is less than a desired precision.

The greedy algorithm written below would require computation of  $\alpha(\mu)$  for  $\#\Xi \times (1 + 2 + \dots + (\#\mathcal{C} - 1))$  values of  $\mu$ . If one choose a large  $\Xi$  set, the cost of the choice of  $\mathcal{C}$  might be too prohibitive. Rather, one could consider a cheap upper bound  $\tilde{\alpha}^{up}(\mu)$  for  $\alpha(\mu)$  and replace  $\alpha(\mu) - \tilde{\alpha}(\mu)$  by the relative surrogate sharpness estimator:  $\frac{\tilde{\alpha}^{up}(\mu) - \tilde{\alpha}(\mu)}{\tilde{\alpha}^{up}(\mu)}$ .

This cheap upper bound could be built using the following: let us define:

$$\mathcal{Y}^{up} = \{y^*(\mu_k), k = 1, \dots, K\}$$

where

$$\{\mu_1, \dots, \mu_K\} = \mathcal{C}$$

and

$$y^*(\mu_k) = \underset{y \in \mathcal{Y}}{\operatorname{arginf}} \left\{ \sum_{q=1}^Q \Theta_q(\mu_k) y_q \right\} \quad (k = 1, \dots, K)$$

We can define:

$$\tilde{\alpha}^{up}(\mu) = \inf_{y \in \mathcal{Y}^{up}} \left\{ \sum_{q=1}^Q \Theta_q(\mu) y_q \right\}$$

which satisfies  $\tilde{\alpha}^{up}(\mu) \geq \alpha(\mu)$  because  $\mathcal{Y}^{up} \subset \mathcal{Y}$ .

Using  $\tilde{\alpha}^{up}(\mu)$  as a surrogate for  $\alpha(\mu)$  reduces fixed complexity of the greedy algorithm to  $K = \#\mathcal{C}$  eigenproblems on  $X$  (for computation of the  $y^*(\mu_k)$ ) and

$\#\Xi$  trivial minimizations per step over reasonably-sized, finite,  $\mathcal{Y}^{up}$ .

## 4 Choice of a reduced basis

In this section, we present two methods for choosing our reduced basis  $\{\zeta_1, \dots, \zeta_N\}$ . These two methods have different offline complexities and lead to different bases. Both nicely fit into our reduced basis framework, and can benefit from the same procedure for error estimation, detailed in section 3.

### 4.1 Greedy procedure

The greedy procedure is treated in [5].

The following algorithm is based on the same heuristic than the one described earlier for choosing  $\mathcal{C}$  in the SCM: at each step of the algorithm, we put one new function in the reduced basis; this function is the solution  $u_{\mu^*}$  of the PDE for a specific value of the parameter choosed so as to give the best expected improvement in the online error bound, i.e. the value of the parameter for which we make the greatest (estimated) error when using "current" reduced basis.

The greedy algorithm for choice of the reduced basis is the following: given  $N \in \mathbb{N}$ , the desired final size of the reduced basis,

- Choose a finite-sized, random, large  $\Xi \subset \mathcal{P}$  sample of parameters.
- Choose  $\mu_1^* \in \mathcal{P}$  at random, and set  $\zeta_1 = u_{\mu_1^*}$ .
- Repeat, for  $n$  from 2 to  $N$ :
  - For each  $\mu \in \Xi$ , compute the online error bound for  $\tilde{u}_\mu$  when using  $\{\zeta_1, \dots, \zeta_{n-1}\}$  as reduced basis.
  - Find  $\mu_n^* \in \Xi$  with the greatest online error bound.
  - Set  $\zeta_n = u_{\mu_n^*}$ .

This procedure requires computation of  $N$  actual solutions, and  $\#\Xi \times N$  online error bounds. It relies crucially on the existence of a cheap online error bound. Of course, instead of fixing a target size for the reduced basis, one can also keep iterating and adding functions to the reduced basis until the greatest online error bound (on the  $\Xi$  sample of parameters) gets below some prescribed value.



## 4.2 Proper Orthogonal Decomposition (POD)-driven procedure

To our knowledge, the following presentation of the Proper Orthogonal Decomposition is new. We refer to [8] for a different presentation, based on singular value decomposition (SVD).

First mode

Suppose, for a soft start, that we want to find a reduced basis containing only one item:  $\{\zeta_1\}$ . Without loss of generality we can always suppose that  $\|\zeta_1\| = 1$ . It seems reasonable to make the choice which makes the least summed squares error when projecting orthogonally the solutions  $u_\mu$  (for all  $\mu \in \mathcal{P}$ ) onto the reduced space  $\text{Span}\{\zeta_1\}$ . Said with formulas, we seek  $\zeta_1$  satisfying:

$$\zeta_1 = \underset{p \in X, \|p\|=1}{\operatorname{arginf}} \int_{\mu \in \Xi} \|u_\mu - \langle u_\mu, p \rangle p\|^2 d\mu$$

The function  $\zeta_1$  will be called our *first POD mode*.

Unfortunately, the integral appearing in the objective function below is, in general, analytically untractable. We can approximate it using a finite sample of parameters  $\Xi \subset \mathcal{P}$ ; the integral then gets replaced by a finite sum over  $\Xi$ :

$$\zeta_1 = \underset{p \in X, \|p\|=1}{\operatorname{arginf}} \sum_{\mu \in \Xi} \|u_\mu - \langle u_\mu, p \rangle p\|^2$$

The set of solutions  $\{u_\mu, \mu \in \Xi\}$  is called the *snapshots* ensemble.

The objective function can then be rewritten:

$$\begin{aligned} \sum_{\mu \in \Xi} \|u_\mu - \langle u_\mu, p \rangle p\|^2 &= \sum_{\mu \in \Xi} \|u_\mu\|^2 - 2 \sum_{\mu \in \Xi} \langle u_\mu, p \rangle^2 + \sum_{\mu \in \Xi} \langle u_\mu, p \rangle^2 \|p\|^2 \\ &= \sum_{\mu \in \Xi} \|u_\mu\|^2 - \sum_{\mu \in \Xi} \langle u_\mu, p \rangle^2 \text{ using that } \|p\|^2 = 1 \end{aligned}$$

And so:

$$\zeta_1 = \underset{p \in X, \|p\|=1}{\operatorname{argmax}} \sum_{\mu \in \Xi} \langle u_\mu, p \rangle^2 \quad (13)$$

The objective function in (13) is a quadratic form in  $p$ , that has to be maximized under a "unity norm" constraint. We know this can be reduced to an eigenproblem.

More specifically, let us introduce a basis  $\{\phi_1, \dots, \phi_{\mathcal{N}}\}$  of  $X$ . As we stressed

in the introduction,  $\mathcal{N}$  is large. We enumerate the snapshots parameters:  $\Xi = \{\mu_1, \dots, \mu_{N_{snap}}\}$ .

Each snapshot  $u(i) = u_{\mu_i}$  ( $i = 1, \dots, N_{snap}$ ) has an expansion in the basis of  $X$ :

$$u(i) = \sum_{j=1}^{\mathcal{N}} u(i)_j \phi_j$$

and we can write down  $M$ , the matrix of the  $u(i)_j$  coefficients:

$$M = \begin{pmatrix} u(1)_1 & u(2)_1 & \dots & u(N_{snap})_1 \\ u(1)_2 & u(2)_2 & \dots & u(N_{snap})_2 \\ \vdots & \vdots & \dots & \vdots \\ u(1)_{\mathcal{N}} & u(2)_{\mathcal{N}} & \dots & u(N_{snap})_{\mathcal{N}} \end{pmatrix}$$

Denote by  $\Omega$ , the (symmetric, positive definite) matrix of the scalar product  $\langle, \rangle$  in  $X$ .

We can write down our objective function as:

$$\sum_{\mu \in \Xi} \langle u_{\mu}, p \rangle^2 = p^T \Omega^T M M^T \Omega p$$

The condition of optimality for (13) is then given by the Lagrange multipliers theorem: there has to exist  $\lambda \in \mathbf{R}$ , so that:

$$\Omega^T M M^T \Omega \zeta_1 = \lambda \Omega \zeta_1 \quad (14)$$

and  $\lambda$  is the optimal value of the problem (as can be seen by left-multiplying (14) by  $\zeta_1^T$ , and using  $\|\zeta_1\|^2 = 1$ ) – which should be made as large as possible. Now left-multiply (14) by  $\Omega^{-T} = \Omega^{-1}$ :

$$M M^T \Omega \zeta_1 = \lambda \zeta_1 \quad (15)$$

so that  $\zeta_1$  is the unit eigenvector associated with the greatest eigenvalue of  $M M^T \Omega$ .

So we have reduced our problem as an eigenproblem of dimension  $\mathcal{N}$ , because  $M M^T \Omega$  is of size  $\mathcal{N} \times \mathcal{N}$ . Since  $\mathcal{N}$  can be large, this could turn out to be highly unpractical. Fortunately, we are to see that we can solve the same problem using eigenproblem of size  $N_{snap} \times N_{snap}$ .

The key is to look for  $\zeta_1$  as a linear combination of *the snapshots*, instead of the

"generic" basis items  $\{\phi_i\}$ .

We write:

$$\zeta_1 = \sum_{i=1}^{N_{snap}} z_{1i} u(i) = Mz_1$$

where  $z_1 = (z_1 1, \dots, z_1 N_{snap})^T$ . Setting  $\zeta_1 = Mz_1$  into (15) yields:

$$MM^T \Omega Mz_1 = \lambda Mz_1$$

and a sufficient condition for this to hold is:

$$M^T \Omega Mz_1 = \lambda z_1$$

so that our new unknown  $z_1 \in \mathbf{R}^{N_{snap}}$  is an eigenvector of  $M^T \Omega M$  (which of size  $N_{snap} \times N_{snap}$ ) associated with its greatest eigenvalue.

The absence of optimality loss when looking for  $\zeta_1$  in  $\text{Range}(M)$  (in other words,  $\zeta_1 = Mz_1$  with  $z_1$  as defined below satisfies (13)) follows clearly from these facts:

1. if  $z$  is a nonzero eigenvector of  $M^T \Omega M$  (associated with  $\lambda \neq 0$ ), then  $\zeta = Mz$  is a nonzero eigenvector of  $MM^T \Omega$  associated with  $\lambda$  ;
2. if  $\zeta$  is a nonzero eigenvector of  $MM^T \Omega$  (associated with  $\lambda \neq 0$ ), then  $z = M^T \Omega \zeta$  is a nonzero eigenvector of  $M^T \Omega M$  associated with  $\lambda$ ;
3.  $M^T \Omega M$  and  $MM^T \Omega$  have the same nonzero eigenvalues, each with the same multiplicity in one and the other.

**Proof:** To see 1., first note that if we had  $Mz = 0$ , then  $M^T \Omega Mz = 0$  and so  $z$  could not be an eigenvector associated with a nonzero eigenvalue of  $M^T \Omega M$ ; thus  $Mz \neq 0$ . Now  $MM^T \Omega (Mz) = M(M^T \Omega Mz) = M\lambda z = \lambda(Mz)$  so that  $Mz$  is an eigenvector of  $MM^T \Omega$  associated with  $\lambda$ . Point 2. is 1. *mutatis mutandis*. Point 3.: it is clear from 1. and 2. that the two matrices have the same nonzero eigenvalues. To get the statement about the multiplicities, suppose that  $\text{Span}\{z_1, \dots, z_k\}$  is the eigenspace of  $M^T \Omega M$  associated with eigenvalue  $\lambda$ . Since  $M^T \Omega M$  is self-adjoint with respect to the Euclidean inner product (i.e., is symmetric), we can suppose that  $\{z_1, \dots, z_k\}$  is orthogonal (for the Euclidean inner product). Now, from 2.,  $\text{Span}\{Mz_1, \dots, Mz_k\}$  is contained in the eigenspace of  $MM^T \Omega$  associated with  $\lambda$ . Moreover, for every  $i \neq j$ :

$$\langle Mz_i, Mz_j \rangle = z_j^T M^T \Omega Mz_i = \lambda z_j^T z_i = 0$$

so that  $\{Mz_1, \dots, Mz_k\}$  is  $\langle, \rangle$ -orthogonal, and therefore linearly independent, and so  $\lambda$  has multiplicity  $k$  in  $MM^T\Omega$ .  $\square$

Second mode

Now let's say we are looking for a set of two functions as our reduced basis:  $\{\zeta_1, \zeta_2\}$ . As earlier, we want  $(\zeta_1, \zeta_2)$  to minimize:

$$\sum_{\mu \in \Xi} \|u_\mu - (\langle u_\mu, \zeta_1 \rangle \zeta_1 + \langle u_\mu, \zeta_2 \rangle \zeta_2)\|^2$$

among the  $(\zeta_1, \zeta_2) \in X \times X$  satisfying:  $\|\zeta_1\| = \|\zeta_2\| = 1$  and  $\langle \zeta_1, \zeta_2 \rangle = 0$ .

In other words, we aim to minimize the overall error (over our sample of parameters  $\Xi$ ) made when orthogonally projecting  $u_\mu$  onto  $\text{Span}\{\zeta_1, \zeta_2\}$ , with  $\{\zeta_1, \zeta_2\}$  an orthonormal family.

Mimicking what we have done earlier, we rewrite our objective function as:

$$\begin{aligned} \sum_{\mu \in \Xi} \|u_\mu - (\langle u_\mu, \zeta_1 \rangle \zeta_1 + \langle u_\mu, p \rangle p)\|^2 &= \sum_{\mu \in \Xi} \|u_\mu\|^2 - 2 \sum_{\mu \in \Xi} \langle u_\mu, \zeta_1 \rangle^2 - 2 \sum_{\mu \in \Xi} \langle u_\mu, p \rangle^2 \\ &\quad + \sum_{\mu \in \Xi} \langle u_\mu, \zeta_1 \rangle^2 \|\zeta_1\|^2 + \sum_{\mu \in \Xi} \langle u_\mu, p \rangle^2 \|p\|^2 \\ &\quad + \sum_{\mu \in \Xi} \langle u_\mu, \zeta_1 \rangle \langle u_\mu, \zeta_2 \rangle \langle \zeta_1, \zeta_2 \rangle \\ &= \text{Constant} - \sum_{\mu \in \Xi} \langle u_\mu, \zeta_1 \rangle^2 - \sum_{\mu \in \Xi} \langle u_\mu, \zeta_2 \rangle^2 \end{aligned}$$

because of the constraints on  $\zeta_1$  and  $\zeta_2$ .

The objective function, apart from the terms that are independent of  $\zeta_1$  and  $\zeta_2$ , is a sum of a function dependent on  $\zeta_1$  only and a function dependent on  $\zeta_2$  only. Thus the joint optimization on  $(\zeta_1, \zeta_2)$  is equivalent to successive optimization on  $\zeta_1$ , then on  $\zeta_2$ . Optimization on  $\zeta_1$  has been done previously. Optimization in  $\zeta_2$  is the following:

$$\zeta_2 = \underset{p \in X, \|p\|=1, \langle p, \zeta_1 \rangle=0}{\operatorname{argmax}} \sum_{\mu \in \Xi} \langle u_\mu, p \rangle^2$$

The objective function has exactly the same matrix expression as the one of the problem defining the first POD mode:

$$\sum_{\mu \in \Xi} \langle u_\mu, p \rangle^2 = p^T \Omega^T M M^T \Omega p$$

The Lagrange multiplier theorem now gives that there exist  $\lambda_2 \in \mathbf{R}$  so that:

$$\Omega^T M M^T \Omega \zeta_2 = \lambda_2 \Omega \zeta_2 \quad (16)$$

As before, left multiplying by  $\zeta_2^T$  gives that  $\lambda_2$  is the optimal value of the problem, which should be made as large as possible.

Left-multiplying (16) by  $\Omega^{-T}$  we get that (16) is equivalent to:

$$M M^T \Omega \zeta_2 = \lambda \zeta_2$$

so  $\lambda_2$  is an eigenvalue of  $M M^T \Omega$ .

Now we have, for every  $x, y \in X$ :

$$\begin{aligned} \langle M M^T \Omega x, y \rangle &= y^T \Omega M M^T \Omega x \\ &= y^T (M M^T \Omega)^T \Omega x \\ &= \langle x, M M^T \Omega y \rangle \end{aligned}$$

so that  $M M^T \Omega$  is self-adjoint with respect to  $\langle, \rangle$ . So we can find a unit eigenvector  $\zeta_2$ , orthogonal to  $\zeta_1$  associated with the second greatest eigenvalue (note that this eigenvalue is equal to  $\lambda_1$  if  $\lambda_1$  has multiplicity  $> 1$ ).

As earlier, we can use the fact that  $M M^T \Omega$  and  $M^T \Omega M$  have the same nonzero eigenvalues, to get  $\zeta_2 = M z_2$  where  $z_2$  is an eigenvector of the smaller matrix  $M^T \Omega M$ . The fact that multiplicities of nonzero eigenvalues in  $M M^T \Omega$  and  $M^T \Omega M$  are the same ensures that this procedure is correct even if  $\lambda_1 = \lambda_2$ .

Following modes, and summary for POD

We can readily repeat what we have done for  $\zeta_2$  to prove:

Let  $\Xi = \{\mu_1, \dots, \mu_{N_{snap}}\}$  be a finite part of  $\mathcal{P}$ . Let  $M$  be the matrix whose columns are the coefficients of  $u_{\mu_1}, \dots, u_{\mu_{N_{snap}}}$  with respect to some basis of  $X$ . Let  $\Omega$  be the matrix of the inner product in  $X$ .

Then, the family  $(\zeta_1, \dots, \zeta_N)$  of vectors in  $X$  that minimizes:

$$\sum_{\mu \in \Xi} \left\| u(\mu) - (\langle u_\mu, \zeta_1 \rangle \zeta_1 + \langle u_\mu, \zeta_2 \rangle \zeta_2 + \dots + \langle u_\mu, \zeta_N \rangle \zeta_N) \right\|^2$$

$$\text{under the constraints } \langle \zeta_i, \zeta_j \rangle = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{else.} \end{cases}$$

is defined by:

$$\begin{cases} \zeta_1 &= \frac{1}{||Mz_1||} Mz_1 \\ &\vdots \\ \zeta_N &= \frac{1}{||Mz_n||} Mz_n \end{cases}$$

where  $z_k$  ( $k = 1, \dots, N$ ) is a nonzero eigenvector, associated with  $k$ -th largest eigenvalue (counting repeated nonsimple eigenvalues) of  $M^T \Omega M$ .

The statement above holds for every  $N \in \mathbf{N}^*$  so that the  $N$ -th largest eigenvalue of  $M^T \Omega M$  is nonzero.

The functions  $(\zeta_1, \dots, \zeta_N)$  defined above are called the first  $N$  POD modes associated with the snapshot ensemble  $\Xi$ .

A note about the (offline) complexity of the POD: for a  $N$ -sized reduced basis, it requires computation of  $N_{snap}$  "expensive" solutions  $u_\mu$  and extraction of the  $N$  leading eigenelements of an  $N_{snap} \times N_{snap}$  symmetric matrix.

## 5 Case of a scalar output

One is often not interested in the full solution  $u_\mu$ , but rather in an "output" calculated from  $u_\mu$ :

$$s_\mu = \eta(u_\mu)$$

where  $\eta : X \rightarrow \mathbf{R}$  is a ( $\mu$ -independent) functional.

We present here some approaches for computation of a surrogate output  $\tilde{s}_\mu$  from  $\tilde{u}_\mu$ .

### 5.1 Primal approach

A trivial approach to the question is to compute  $\tilde{u}_\mu$ , and then, online:

$$\tilde{s}_\mu = \eta(\tilde{u}_\mu) \tag{17}$$

Done naively, this procedure might be  $\dim X$ -dependent, since evaluation of  $\eta$  may require access to every component of  $\tilde{u}_\mu$  in the original  $X$  space.

A way of overcoming this problem, if  $\eta$  is a linear functional, is to write that:

$$\tilde{s}_\mu = \eta \left( \sum_{n=1}^N c_n(\mu) \zeta_n \right) = \sum_{n=1}^N c_n(\mu) \eta(\zeta_n)$$

and the offline/online procedure becomes:

- offline: compute and store  $\eta(\zeta_n)$ ,  $n = 1, \dots, N$
- online: find  $c_n(\mu)$  ( $n = 1, \dots, N$ ) satisfying (2) using procedure described in 2.2 and compute

$$\tilde{s}_\mu = \sum_{n=1}^N c_n(\mu) \eta(\zeta_n) \quad (18)$$

And now the online phase is still  $\dim X$ -independent.

For a quadratic output functional, one proceeds the same way using:

$$\tilde{s}_\mu = \sum_{n=1}^N \sum_{n'=1}^N c_n(\mu) c_{n'}(\mu) \eta_b(\zeta_n, \zeta_{n'})$$

where  $\eta_b$  is the symmetric bilinear form associated with  $\eta$ .

This time, one has to compute and store (offline)  $\eta_b(\zeta_n, \zeta_{n'})$ , for  $n, n' = 1, \dots, N$ .

The same procedure can easily be generalized for polynomial outputs.

Error bound for the output

A straightforward majoration for  $|s_\mu - \tilde{s}_\mu|$  can be found under a Lipschitz hypothesis for  $\eta$ : let  $L > 0$  be so that, whatever  $v, w \in X$ :

$$|\eta(v) - \eta(w)| \leq L \|v - w\|$$

(for example, if  $\eta$  is a continuous linear functional,  $L$  can be chosen as the dual norm of  $\eta$ ).

Then we have:

$$|s_\mu - \tilde{s}_\mu| \leq L \varepsilon_\mu \quad (19)$$

## 5.2 Primal/dual approach

For linear outputs, error bound (19) can be improved, at the expense of, in some cases, an increase in computation time and program complexity. The idea, present in [5], is to introduce the *dual* problem of (1) associated with  $\eta$ :

$$a_\mu(v, \psi_\mu) = -\eta(v) \quad \forall v \in X \quad (20)$$

whose unknown is  $\psi_\mu \in X$ .

This dual problem is an affine parametric PDE which we can apply the RB framework on. That is, given a reduced basis  $\{\zeta_1^{du}, \dots, \zeta_{N_{du}}^{du}\}$  and a reduced space  $\tilde{X}^{du} = \text{Span}\{\zeta_1^{du}, \dots, \zeta_{N_{du}}^{du}\}$ , find an approximation  $\tilde{\psi}_\mu \in \tilde{X}^{du}$  to  $\psi_\mu$  satisfying:

$$a_\mu(v, \tilde{\psi}_\mu) = -\eta(v) \quad \forall v \in \tilde{X}^{du}$$

with error bound:

$$\|\tilde{\psi}_\mu - \psi_\mu\| \leq \varepsilon_\mu^{du}$$

After having found reduced basis solutions for the primal and dual problems, the following surrogate output can be taken:

$$\tilde{s}_\mu = \eta(\tilde{u}_\mu) - r_\mu(\tilde{\psi}_\mu) \quad (21)$$

The error bound on the output now reads:

$$\|\tilde{s}_\mu - s_\mu\| \leq \alpha(\mu) \varepsilon_\mu \varepsilon_\mu^{du} \quad (22)$$

Using (21) instead of (17) gives us the quadratic error bound (22) instead of the linear (19) for only (roughly) twice the price in computation and storage requirements. And if the problem is self-dual (i.e.  $a_\mu$  is symmetric and  $\eta = -f$ ), we can take  $N = N_{du}$ ,  $\zeta_i = \zeta_i^{du} \quad \forall i = 1, \dots, N$  and get this "squared" effect for free, by just adding up the corrective term in (21).

**Proof of (22):** We omit the  $\mu$  subscripts.

$$\begin{aligned} |s - \tilde{s}| &= |\eta(u - \tilde{u}) - a(\tilde{u}, \tilde{\psi}) + f(\tilde{\psi})| \\ &= |-a(u - \tilde{u}, \psi) - a(\tilde{u}, \tilde{\psi}) + f(\tilde{\psi})| \text{ from (20)} \\ &= |-a(u, \psi) + f(\tilde{\psi}) + a(\tilde{u}, \psi - \tilde{\psi})| \text{ by reordering terms} \\ &= |f(\tilde{\psi} - \psi) + a(\tilde{u}, \psi - \tilde{\psi})| \text{ from (1)} \\ &= |r(\psi - \tilde{\psi})| \text{ by the very definition of } r \\ &\leq \|r\|_{X'} \|\psi - \tilde{\psi}\| \leq \|r\|_{X'} \varepsilon^{du} = \alpha \varepsilon \varepsilon^{du} \text{ from the definition of } \varepsilon. \quad \square \end{aligned}$$

We shall now pass to the offline/online procedure: we first require computation of  $r_\mu(\tilde{\psi}_\mu)$  in (21).

Expanding  $\tilde{\psi}_\mu$  on the dual reduced basis:

$$\tilde{\psi}_\mu = \sum_{p=1}^{N_{du}} q_p(\mu) \zeta_p^{du}$$



we have:

$$r_\mu(\tilde{\psi}_\mu) = \sum_{p=1}^{N_{du}} q_p(\mu) r_\mu(\zeta_p^{du}) \quad (23)$$

and, for each  $p = 1, \dots, N_{du}$ , we have, from definition of  $r_\mu$ , affine expansion of  $a_\mu$  and reduced basis expansion (2) of  $\tilde{u}_\mu$ :

$$r_\mu(\zeta_p^{du}) = f(\zeta_p^{du}) - \sum_{q=1}^Q \Theta_q(\mu) \sum_{n=1}^N c_n(\mu) a_q(\zeta_n, \zeta_p^{du}) \quad (24)$$

So our full offline/online procedure for primal/dual evaluation and error bound on  $s_\mu$  reads:

- offline:
  - perform offline-stage RB approximation (part 2.2) and error estimation (part 3) for the two problems (1) and (20) ;
  - compute and store  $\eta(\zeta_n)$  for  $n = 1, \dots, N$  ;
  - compute and store  $f(\zeta_p^{du})$  ( $p = 1, \dots, N_{du}$ ) and  $a_q(\zeta_n, \zeta_p^{du})$  ( $q = 1, \dots, Q; n = 1, \dots, N; p = 1, \dots, N_{du}$ ) ;
- online:
  - perform online-stage RB approximation and error estimation for the two problems (1) and (20); this gives respective expansions  $\{c_n(\mu)\}_{n=1, \dots, N}$  and  $\{q_p(\mu)\}_{p=1, \dots, N_{du}}$  of  $\tilde{u}_\mu$  and  $\tilde{\psi}_\mu$  in, respectively primal reduced basis  $\{\zeta_1, \dots, \zeta_N\}$  and dual reduced basis  $\{\zeta_1^{du}, \dots, \zeta_{N_{du}}^{du}\}$ , together with  $\alpha(\mu)$ ,  $\varepsilon_\mu$  (error bound for primal problem) and  $\varepsilon_\mu$  (error bound for dual problem) ;
  - compute  $\eta(\tilde{\mu})$  as left-hand side of (18) ;
  - compute  $r_\mu(\tilde{\psi}_\mu)$  using (23) and (24) ;
  - output  $\tilde{s}_\mu$  using (21) ;
  - give error bound (22).

### 5.3 Output-aware choice of basis: goal-oriented POD

When one is interested in a special output of  $u_\mu$ , it could be wise to choose an *adapted* reduced basis for this output. For instance, if the output depends strongly on the values of  $u_\mu(x)$  for  $x$  in some region  $R_1$  contained in the domain

of calculus, and much less strongly for  $x$  in some region  $R_2$ , then a well-chosen reduced basis should lead to a precise  $\tilde{u}_\mu$  approximation in  $R_1$ , which could be much less precise in  $R_2$ .

The goal-oriented POD method, presented in [1] in the context of linear-time-invariant (LTI) dynamical systems and balanced truncation, is a variation on POD taking the output functional in account.

Goal-oriented POD proceeds by finding  $(\zeta_1, \dots, \zeta_N) \in X$  minimizing the sum of squared differences (for different values of a parameter in a finite sample  $\#\Xi$ ) between the actual output and the output evaluated on the RB approximation taking  $\{\zeta_1, \dots, \zeta_N\}$  so as to minimize

$$\inf_{\mu \in \Xi} \sum_{\mu \in \Xi} \|\eta(u_\mu) - \eta(\tilde{u}_\mu)\|^2 \quad (25)$$

under the constraints:

$$\left\{ \begin{array}{l} \tilde{u}_\mu = \sum_{n=1}^N c_n(\mu) \zeta_n \quad \forall \mu \in \Xi \\ \sum_{n=1}^N \left[ \sum_{q=1}^Q \Theta_q(\mu) a_q(\zeta_n, \zeta_p) \right] c_n(\mu) = f(\zeta_p) \quad \forall \mu \in \Xi, \forall p = 1, \dots, N \\ \langle \zeta_n, \zeta_m \rangle = \begin{cases} 1 & \text{if } n = m \\ 0 & \text{else} \end{cases} \quad \forall n, m = 1, \dots, N \end{array} \right.$$

We believe it is possible to adapt the optimization framework described in [1] to solve (a weaker variant) of the above problem for  $(\zeta_1, \dots, \zeta_N)$ . However this optimization problem seems rather difficult (non-convex, with local non-global minima), thus requiring much more computational burden than classic POD.

## Conclusion

We have presented a PDE-based, deterministic approach for model reduction called reduced basis approximation. The advantages of this approach are: a certified, computationally-efficient bound for reduction error, and a wide flexibility in the reduced basis chosen so as to enable output-aware choices of basis. Its main drawback is the need to make assumptions on the model that is being reduced. However work have been done to extend the reduced basis, and its associate procedures (error estimation, basis selection) to non-affinely

parametrized PDEs ([2]) or non-linear, time-dependent PDEs (see [4] for instance).

Future work should include further extensions to other PDEs with different parameters, and more output-aware basis choice procedures.

## References

- [1] T. Bui-Thanh, K. Willcox, O. Ghattas, and B. van Bloemen Waanders. Goal-oriented, model-constrained optimization for reduction of large-scale systems. *Journal of Computational Physics*, 224(2):880–896, 2007.
- [2] M.A. Grepl, Y. Maday, N.C. Nguyen, and A.T. Patera. Efficient reduced-basis treatment of nonaffine and nonlinear partial differential equations. *Mathematical Modelling and Numerical Analysis*, 41(3):575–605, 2007.
- [3] D.B.P. Huynh, G. Rozza, S. Sen, and A.T. Patera. A successive constraint linear optimization method for lower bounds of parametric coercivity and inf-sup stability constants. *Comptes Rendus Mathematique*, 345(8):473–478, 2007.
- [4] N.C. Nguyen, G. Rozza, and A.T. Patera. Reduced basis approximation and a posteriori error estimation for the time-dependent viscous Burgers’ equation. *Calcolo*, 46(3):157–185, 2009.
- [5] N.C. Nguyen, K. Veroy, and A.T. Patera. Certified real-time solution of parametrized partial differential equations. *Handbook of Materials Modeling*, pages 1523–1558, 2005.
- [6] J. Nocedal and S.J. Wright. *Numerical optimization*. Springer Verlag, 1999.
- [7] A.M. Quarteroni and A. Valli. *Numerical approximation of partial differential equations*. Springer, 2008.
- [8] S. Volkwein. Proper orthogonal decomposition and singular value decomposition. 1999.