

## 1 REDUCED ORDER MODEL: HEAT EQUATION

We present now the formulation for the Reduced Order Model (ROM) of our FE problem. Our aim is to be able to, for any unseen parameter value, assemble and solve a smaller problem in terms of the involved algebraic operators dimensions. In order to do so, first we need to obtain certain algebraic structures which capture the essence of the problem at hand. We do so by sampling the original problem at certain parameter values and processing snapshots of the solution and the operators, obtained from the solution of the FOM problem. This is called the *offline phase*.

Later on, we exploit these static structures to build reduced operators which capture the dynamics of the problem sufficiently well, solve a smaller algebraic system and then recover the solution in the original mesh variables, in order to postprocess it. This is called the *online phase*.

We will use the Reduced Basis Method (RB) to construct an ad-hoc problem-based basis to represent our ROM, and the Discrete Empirical Interpolation Method (DEIM) and its matrix version (MDEIM) to build suitable approximations of the algebraic operators involved.

The continuous reduced problem formulation is skipped since it will not be used. Therefore, we jump directly into the discrete problem. We recall that we are focused at reducing the problem for the homogeneous component  $\hat{u}(x)$  of our solution, i.e., the weak formulation given by the system of equations (??).

### 1.1 Reduced Basis Method: A Naive Galerkin Approach

We have included in the title the word *naive* because we are going to define the reduction problem in a very blunt way, where many inefficiencies will show up. We do so to motivate the additional reduction procedures we will include later on, specially regarding the operators.

In the reduced context, we use a finite space  $V_N \subset V_h$ , where we can represent the solution as the linear combination of a set of orthonormal<sup>1</sup> problem-based basis functions  $\psi_i(x)$ ,

$$\hat{u}_N(x) = \sum_j^N \hat{u}_{Nj} \psi_j(x). \quad (1)$$

These basis functions  $\psi_i(x)$  have global support, since their goal is to capture the problem dynamics. This is why we usually expect or desire to have  $N \ll N_h$ , since we want to reduce the number of basis functions we need to represent our solution.

To maintain focus and in favor of generality, let us assume at this point that the global basis functions are given, and that they are zero at the boundary. We will give details on how to obtain them later on.

#### 1.1.1 Reduced Space Projection

Since we can represent any function in terms of our FE basis functions  $\varphi_i(x)$ , we have a linear mapping between the problem-based functions and the nodal basis functions. This allows us to establish the following relation between the problem solution in the reduced and original spaces,

$$\hat{u}_h = \mathbb{V} \hat{u}_N. \quad (2)$$

1. If they were not, we can always make them so via Gram-Schmidt.

The entries of the  $\mathbb{V}$  matrix represent the coefficients of the global basis representation in the FE basis.

$$\psi_i(x) = \sum_j [\mathbb{V}]_{ji} \varphi_j(x) \quad (3)$$

#### 1.1.2 Discrete Reduced Problem Assembly

In theory, to assemble the reduced problem operators, one could actually compute the inner products defined in Equations (??) with these new basis functions  $\psi_i(x)$ . In practice, it is more convenient to project the algebraic FOM operators with matrix-matrix and matrix-vector products into the reduced space,

$$\mathbf{M}_N^{n+1} = \mathbb{V}^T \mathbf{M}_h^{n+1} \mathbb{V}, \quad (4a)$$

$$\mathbf{A}_N^{n+1} = \mathbb{V}^T \mathbf{A}_h^{n+1} \mathbb{V}, \quad (4b)$$

$$\mathbf{F}_N^{n+1} = \mathbb{V}^T \mathbf{F}_h^{n+1}, \quad (4c)$$

$$\mathbf{F}_{g,N}^{n+1} = \mathbb{V}^T \mathbf{F}_{g,h}^{n+1}, \quad (4d)$$

since the assembly of matrices based on the FE basis is easily done in parallel due to their local support, whereas the integration of functions with global support is not necessarily computationally efficient.

By doing so, we find the following ROM for the time evolution problem,

$$\mathbf{M}_N \hat{\mathbf{u}}_N^{n+1} + \Delta t \mathbf{A}_N \hat{\mathbf{u}}_N^{n+1} = \mathbf{F}_{\hat{\mathbf{u}}_N}^n \quad (5a)$$

$$+ \Delta t \mathbf{F}_N^{n+1} + \Delta t \mathbf{F}_{g,N}^{n+1},$$

$$\hat{\mathbf{u}}_N^0 = \hat{\mathbf{u}}_{N,0}. \quad (5b)$$

If we collect terms and factor out the unknowns we get a linear system, this time in the reduced space, to be solved for each timestep to advance the solution,

$$\mathbf{K}_N^{n+1} \hat{\mathbf{u}}_N^{n+1} = \mathbf{b}_N^{n+1}, \quad (6a)$$

$$\mathbf{K}_N^{n+1} = \mathbf{M}_N^{n+1} + \Delta t \mathbf{A}_N^{n+1}, \quad (6b)$$

$$\mathbf{b}_N^{n+1} = \mathbf{F}_{\hat{\mathbf{u}}_N}^n + \Delta t \mathbf{F}_N^{n+1} + \Delta t \mathbf{F}_{g,N}^{n+1}, \quad (6c)$$

$$\hat{\mathbf{u}}_N^0 = \hat{\mathbf{u}}_{N,0}. \quad (6d)$$

#### Time-Discretization Forcing Term

The forcing term  $\mathbf{F}_{\hat{\mathbf{u}}_h}^n$  due to the time discretization has been intentionally left out in the previous section. We could naively project it too,

$$\mathbf{F}_{\hat{\mathbf{u}}_N}^n = \mathbb{V}^T \mathbf{F}_{\hat{\mathbf{u}}_h}^n, \quad (7)$$

but this would force us to reconstruct the FE vector of the ROM at each time-step, making the integration cumbersome. To go around this issue, we can exploit the algebraic representation of  $\mathbf{F}_{\hat{\mathbf{u}}_h}^n$ , given in Equation (??) in terms of the mass matrix. Consequently, the forcing term due to the time discretization takes the following form in the ROM:

$$\mathbf{F}_{\hat{\mathbf{u}}_N}^n = \begin{cases} \mathbf{M}_N^{n+1} \hat{\mathbf{u}}_N^n, & \text{BDF-1} \\ \frac{1}{2} \mathbf{M}_N^{n+1} \hat{\mathbf{u}}_N^n - \frac{3}{4} \mathbf{M}_N^{n+1} \hat{\mathbf{u}}_N^{n-1}, & \text{BDF-2} \end{cases} \quad (8)$$

#### 1.1.3 Boundary and Initial Conditions

The computation of the initial condition  $\hat{\mathbf{u}}_{N,0}$  is to be done in two steps. First, the initial condition  $\hat{\mathbf{u}}_{h,0}$  in the FOM space needs to be computed as a FE vector by means of the techniques explained in Section ?? . Then, to obtain  $\hat{\mathbf{u}}_{N,0}$ , the FE representation  $\hat{\mathbf{u}}_{h,0}$  needs to be projected unto the reduced

space. Since we are dealing with an orthonormal basis, we can use the matrix  $\mathbb{V}$  to project the FE vector departing from the FOM-ROM relation given in Equation (2),

$$\mathbb{V}^T \hat{\mathbf{u}}_{h,0} = \underbrace{\mathbb{V}^T \mathbb{V}}_{\mathbb{I}} \hat{\mathbf{u}}_{N,0} \rightarrow \hat{\mathbf{u}}_{N,0} = \mathbb{V}^T \hat{\mathbf{u}}_{h,0}. \quad (9)$$

Regarding the spatial boundary conditions, at this point of the naive reduction scheme, two facts come into play, which we first present and then put together:

- 1) The weak form has homogeneous boundary conditions.
- 2) The ad-hoc basis  $\psi_i(x)$  take value zero at the boundaries.

These two facts imply that the projection of the operators, system of Equations (4), does not break the original constraint of homogeneous boundary conditions; and that the linear expansion of the solution  $\hat{u}_N(x)$  in the span of  $V_N$ , Equation (1), is always true.

Once the reduced homogeneous solution  $\hat{u}_N(x)$  is obtained, it can be brought back to the original space  $V_h$  via Equation (2), and then add on top of it the FE representation of the Dirichlet lifting, to obtain the final solution,

$$u(x, t; \mu) \simeq u_h(t, \mu) = \mathbb{V} \hat{u}_N(t, \mu) + g_h(t, \mu). \quad (10)$$

#### 1.1.4 Conclusions

At this point, the naive reduction scheme for the RB-ROM has been defined and explained. However, some aspects of it remain unattended:

- The construction of the problem-based basis functions  $\psi_i(x)$ . There are many methods to build them, and which combination of them we chose defines which reduction method we are applying, along with their advantages and requirements.
- The *offline-online decomposition*: to uncouple the usage of FOM operators in as much as possible from the integration of the ROM. Ideally, no FOM structures should be assembled during the online phase. In this naive scheme, we are clearly not meeting such requirement, since we need to assemble all the FOM operators and then project them into  $V_N$  for each timestep.

Thus, some additional sections need to be brought up, in order to complete our definition of the reducing scheme. We now treat the construction of the basis functions, Section 1.2; and the approximation of the algebraic operators, Section 1.3.

## 1.2 Reduced Basis Construction

Hereby we layout the details of the basis construction, that is, the definition of the  $\psi_i(x)$  functions. There are many techniques to build such basis, but since we are laying out a toy problem, we explain an automatic and out-of-the-box technique: the nested POD approach.

On paper, the most simple basis anyone could come up with is a collection of solution snapshots for different parameter values and timesteps,

$$\begin{aligned} \Psi_{\hat{u}} &:= [\psi_i] \\ &= [\hat{u}_h(t^0; \mu_0), \hat{u}_h(t^1; \mu_0), \dots, \hat{u}_h(t^j; \mu_i), \dots], \\ &= [\Psi_{\hat{u}}(\mu_0), \Psi_{\hat{u}}(\mu_1), \dots, \Psi_{\hat{u}}(\mu_{N_\mu})] \end{aligned} \quad (11)$$

Yet, this basis  $\Psi_{\hat{u}}$  is unpractical from a computational point of view: we could not possibly store all the basis vectors and neither compute efficiently all the required algebraic operations with them. Additionally, it would probably lead to ill-posed linear systems, since the vectors are almost linearly dependent.

However, since all these vectors arise from the same PDE (although for different parametrizations of it), and for each timestep the solution is close to the previous one, in a way, we could expect there to be a lot of repeated information inside each  $\Psi_{\hat{u}}(\mu_i)$ , and consequently, inside  $\Psi_{\hat{u}}$ . We can exploit this fact by using a compression algorithm, such as the Proper Orthogonal Decomposition, to find a set of vectors which capture sufficiently good the span of  $\Psi_{\hat{u}}$ , and yet do so with less number of basis vectors.

We call this the *method of snapshots*.

### 1.2.1 POD Space Reduction

This section is left intentionally short for the moment. Many high quality references exist to explain what is the Proper Orthogonal Decomposition (POD), why does it work and which limitations does it have.

We can see it as an enhanced Gram-Schmidt orthonormalization procedure: not only an orthonormal basis is obtained, but additionally the output basis vectors recover hierarchically the span of the original space given by the input vectors.

Let us define the  $POD : X \rightarrow Y$  function between two normed spaces, such that  $Y \subseteq X$ , which takes as inputs a collection of  $N_S$  vectors and a prescribed tolerance error  $\varepsilon_{POD}$ ,

$$[\psi_i]_{i=1}^{N_{POD}} = POD \left( [\varphi_i]_{i=1}^{N_S}, \varepsilon_{POD} \right). \quad (12)$$

with  $\varphi_i \in X$  and  $\psi_i \in Y$ . This function returns a collection of orthonormal  $N_{POD}$  vectors, whose span is a subset of the input vector span.

Internally, the POD is solving an optimality approximation problem in the  $L^2$  norm. Therefore, with a slight notation abuse, one could conceptually say

$$\text{span}(\varphi_i) = \text{span}(\psi_i) + \varepsilon_{POD}, \quad (13)$$

that is, the representation of any vector in the original space can be reconstructed to a point with the output POD basis, and the error in the  $L^2$  norm should be less or equal to  $\varepsilon_{POD}$ .

There is a relation between the prescribed tolerance error  $\varepsilon_{POD}$  and the outcoming number of basis elements  $N_{POD}$ ,

$$N_{POD} = N_{POD}(\varepsilon_{POD}). \quad (14)$$

This functional relationship usually shows exponential decay, or a sharp drop beyond a given number of elements. That is, if a problem is reduceable, beyond a given number of elements, adding more will not improve significantly the approximation bondness.

Alternatively to a prescribed error, one could directly ask for a prescribed number of basis functions  $N_{POD}^* \leq N_S$ ,

$$[\psi_i]_{i=1}^{N_{POD}^*} = POD \left( [\varphi_i]_{i=1}^{N_S}, N_{POD}^* \right). \quad (15)$$

### 1.2.2 Nested POD Basis Construction

A simple procedure to leverage the compression properties of the POD function, and the nested structure of our PDE with respect to time and the parameters; is to first build certain POD basis from the snapshots of the solution at different timesteps for a fixed parameter value,

$$\begin{aligned}\Psi_{\mu_0} &= \text{POD}_\varepsilon \left( \left[ \hat{u}_h(t^0; \mu_0), \dots, \hat{u}_h(t^T; \mu_0) \right] \right), \\ \Psi_{\mu_1} &= \text{POD}_\varepsilon \left( \left[ \hat{u}_h(t^0; \mu_1), \dots, \hat{u}_h(t^T; \mu_1) \right] \right), \\ &\dots, \\ \Psi_{\mu_{N_\mu}} &= \text{POD}_\varepsilon \left( \left[ \hat{u}_h(t^0; \mu_{N_\mu}), \dots, \hat{u}_h(t^T; \mu_{N_\mu}) \right] \right).\end{aligned}$$

Then, all the  $\mu$ -fixed POD basis, which sum up the information contained in the time evolution direction, are compressed again using a POD,

$$\Psi := \Psi = \text{POD}_\varepsilon \left( \left[ \Psi_{\mu_0}, \Psi_{\mu_1}, \dots, \Psi_{\mu_{N_\mu}} \right] \right).$$

In the end, this gives us a unique basis  $\Psi = [\psi_i] = \mathbb{V}$ , which contains information for parameter variations and time evolution. And above all, it has been built in a computationally efficient manner, at least in terms of storage.

Different error tolerances could be prescribed at the time and parameter compression stages.

We say this to be an automatic and out-of-the-box procedure because it does not require further complications beyond the storage of the snapshots and the implementation of the POD algorithm. It only demands the definition of a collection of parameter values to solve for. This can be done with random sampling techniques, or if some physically-based knowledge is available, a custom selection of the parameters for ranges in which we know the solution will strongly vary.

Naturally, more involved procedures exist to create the final basis  $\Psi$ , but they demand the capacity to evaluate during the creation of the basis how good or bad the new basis is performing at capturing the problem dynamics, or the determination of sharp *a priori* error estimators.

We leave this for later.

### 1.3 System Approximation

Regarding the assembly of the reduced operators, if already during the offline stage, where the FOM problem is tackled, we had to assemble all the discrete operators for each timestep; during the online stage we have to additionally project them unto the reduced space too, as shown in Equations (4), increasing the overall cost of the integration procedure.

This will be our main motivation to include a system approximation procedure, with the goal of speeding up the construction of the operators.

We talk about *parameter and time separable* problems, or the existence of an *affine decomposition*, when the spatial operators (bilinear or linear forms), which depend on time and parameter values, present the following functional form,

$$A_h(t, \mu) = \sum_q^{Q_a} \Theta_q^a(t, \mu) A_{h,q}, \quad (16)$$

where the coefficient functions are real-valued,  $\Theta_q^a(t, \mu) \in \mathbb{R}$ , and the operator basis elements  $A_{h,q}$  are parameter-independent.

This expansion can be used for both matrices or vectors provided the topology of the mesh does not change in time. If this is the case, the matrices can be transformed into vectors, and later on brought back to matrix form once any necessary operation has been carried out.

If we had such a decomposition, once we had computed the basis matrix  $\mathbb{V}$ , we could project each element of the operator basis  $A_{h,q}$  to obtain an expression for the reduced operator,

$$\begin{aligned}A_N(t, \mu) &= \mathbb{V}^T A_h(t, \mu) \mathbb{V} \\ &= \sum_q^{Q_a} \Theta_q^a(t, \mu) \mathbb{V}^T A_{h,q} \mathbb{V} \\ A_N(t, \mu) &= \sum_q^{Q_a} \Theta_q^a(t, \mu) A_{N,q}.\end{aligned} \quad (17)$$

Since  $A_{N,q}$  is fixed, provided that we had a way to evaluate each  $\Theta_q^a(t, \mu)$ , we would be able to build the reduced operator for a given parameter for each timestep without having to use any FOM operator.

#### 1.3.1 Discrete Empirical Interpolation Method

Naturally, not many problems are likely to present a separable form as the one shown above. Even our simple linear heat equation problem, due to the time-deformation of the mesh, cannot be presented in such a form.

To tackle this issue, we use the Discrete Empirical Interpolation Method (DEIM). This method is a numerical extension of its analytical sibling, the Empirical Interpolation Method (EIM). Basically, it mimicks the idea of creating a basis for the solution space, but this time centered around the operator space. By means of a nested POD as we explained in Section 1.2.2, if we replace the solution snapshots with operator snapshots, we can build the static and problem-dependent basis  $A_{h,q}$ .

Since we will be creating the operator basis with an approximation technique, an error is expected in the reconstruction of the actual operator, and so we introduce the notation  $A_h^m(t, \mu)$  to reference the approximation of the operator via the (M)DEIM algorithm,

$$A_h(t, \mu) \simeq A_h^m(t, \mu) = \sum_q^{Q_a} \Theta_q^a(t, \mu) A_{h,q}. \quad (18)$$

Naturally, this idea leads to the concept of approximated reduced operators,

$$A_N(t, \mu) \simeq A_N^m(t, \mu) = \sum_q^{Q_a} \Theta_q^a(t, \mu) A_{N,q}, \quad (19)$$

which is the approximation of the reduced operators when the basis comes from an approximation of the operator snapshots.

#### 1.3.2 Discussion on the Approximation Target

There are reasons to approximate  $A_h$  instead of directly  $A_N$ , but I still need to wrap my head around them.

### 1.3.3 Evaluation of the Coefficient Functions

To evaluate the  $\Theta_q^a(t, \mu)$  functions, we set and solve an interpolation problem; that is, we enforce the approximation to actually match certain elements of the operator,

$$[A_h(t, \mu)]_k = [A_h^m(t, \mu)]_k = \sum_q^{Q_a} \Theta_q^a(t, \mu) [A_{h,q}]_k \quad (20)$$

for certain indices  $k \in \mathcal{I}_a$  within a collection of  $Q_a$  selected indices  $\mathcal{I}_a$ . The notation  $[A_h(t, \mu)]_k$  stands for the value of the operator at the given mesh node  $k$ . In the FE context it can be obtained by integrating the weak form locally.

This leads to a system with the same number of unknowns as equations, where each  $\Theta_q^a(t, \mu)$  is unknown. The indices  $\mathcal{I}_a$  are selected during the offline stage, according to error reduction arguments of the separable form (18),

$$e_a(t, \mu) = \|A_h(t, \mu) - A_h^m(t, \mu)\|. \quad (21)$$

### 1.4 Hyper Reduced Basis Method

With an approximation of the reduced operators available, we can define yet another algebraic problem to integrate and obtain the reduced solution in time for a given parameter value,

$$\mathbf{M}_N^{m,n+1} \hat{\mathbf{u}}_N^{n+1} + \Delta t \mathbf{A}_N^{m,n+1} \hat{\mathbf{u}}_N^{n+1} \quad (22a)$$

$$= \mathbf{F}_{\hat{\mathbf{u}}_N}^n + \Delta t \mathbf{F}_N^{m,n+1} + \Delta t \mathbf{F}_{g,N}^{m,n+1},$$

$$\hat{\mathbf{u}}_N^0 = \hat{\mathbf{u}}_{N,0}. \quad (22b)$$

If we collect terms and factor out the unknowns we get a linear system, in the reduced space and with approximated operators, to be solved for each timestep to advance the solution,

$$\mathbf{K}_N^{m,n+1} \hat{\mathbf{u}}_N^{n+1} = \mathbf{b}_N^{m,n+1}, \quad (23a)$$

$$\mathbf{K}_N^{m,n+1} = \mathbf{M}_N^{m,n+1} + \Delta t \mathbf{A}_N^{m,n+1}, \quad (23b)$$

$$\mathbf{b}_N^{m,n+1} = \mathbf{F}_{\hat{\mathbf{u}}_N}^n + \Delta t \mathbf{F}_N^{m,n+1} + \Delta t \mathbf{F}_{g,N}^{m,n+1}, \quad (23c)$$

$$\hat{\mathbf{u}}_N^0 = \hat{\mathbf{u}}_{N,0}. \quad (23d)$$

Each of the operators present in the problem,  $\mathbf{M}_N$ ,  $\mathbf{A}_N$ ,  $\mathbf{F}_N$  and  $\mathbf{F}_{g,N}$  will have associated an operator basis and will require the solution of the interpolation problem (20) to be solved for each timestep and parameter value. Although this could still seem like a costly procedure, if the operators are actually reduceable, the number of basis functions  $Q_m, Q_a, Q_f, Q_{f,g}$  should be small, and thus way simpler problems than assembling the whole operator and the carrying out the projection.

Once the reduced homogeneous solution  $\hat{u}_N^m(x)$  is obtained with approximated operators, it can be brought back to the original space  $V_h$  via Equation (2), and then add on top of it the FE representation of the Dirichlet lifting, to obtain the final solution,

$$u(x, t; \mu) \simeq u_h^m(t, \mu) = \mathbb{V} \hat{u}_N^m(t, \mu) + g_h(t, \mu). \quad (24)$$

### 1.5 Certification of the Reduction Procedure

Should this have some dedicated paragraphs at the beginning? I think it should.

TBD.

- Computation of error bounds between FOM and ROM without actually computing the FOM.