

Hessian-Based Model Reduction for Large-Scale Data Assimilation Problems*

Omar Bashir¹, Omar Ghattas², Judith Hill³,
Bart van Bloemen Waanders³, and Karen Willcox¹

¹ Massachusetts Institute of Technology, Cambridge MA 02139, USA
`bashir@mit.edu, kwillcox@mit.edu`

² The University of Texas at Austin, Austin TX 78712
`omar@ices.utexas.edu`

³ Sandia National Laboratories, Albuquerque NM 87185
`jhill@sandia.gov, bartv@sandia.gov`

Abstract. Assimilation of spatially- and temporally-distributed state observations into simulations of dynamical systems stemming from discretized PDEs leads to inverse problems with high-dimensional control spaces in the form of discretized initial conditions. Solution of such inverse problems in “real-time” is often intractable. This motivates the construction of reduced-order models that can be used as surrogates of the high-fidelity simulations during inverse solution. For the surrogates to be useful, they must be able to approximate the observable quantities over a wide range of initial conditions. Construction of the reduced models entails sampling the initial condition space to generate an appropriate training set, which is an intractable proposition for high dimensional initial condition spaces unless the problem structure can be exploited. Here, we present a method that extracts the dominant spectrum of the input-output map (i.e. the Hessian of the least squares optimization problem) at low cost, and uses the principal eigenvectors as sample points. We demonstrate the efficacy of the reduction methodology on a large-scale contaminant transport problem.

Keywords: Model reduction; data assimilation; inverse problem; Hessian matrix; optimization.

1 Introduction

One important component of Dynamic Data Driven Application Systems (DDDAS) is the continuous assimilation of sensor data into an ongoing simulation. This inverse problem can be formulated as an optimal control problem,

* Partially supported by the National Science Foundation under DDDAS grants CNS-0540372 and CNS-0540186, the Air Force Office of Scientific Research, and the Computer Science Research Institute at Sandia National Laboratories. Sandia is a multiprogram laboratory operated by Sandia Corporation, a Lockheed-Martin Company, for the US Department of Energy under Contract DE-AC04-94AL85000.

in which the controls are the initial conditions, the constraints are the state equations describing the dynamics of the system, and the objective is the difference between the state observations and those predicted by the state equations, measured in some appropriate norm.

When the physical system being simulated is governed by partial differential equations in three spatial dimensions and time, the forward problem alone (i.e. solution of the PDEs for a given initial condition) may require many hours of supercomputer time. The inverse problem, which requires repeated solution of the forward problem, may then be out of reach in situations where rapid assimilation of the data is required. In particular, when the simulation is used as a basis for forecasting or decision-making, a *reduced model* that can execute much more rapidly than the high-fidelity PDE simulation is then needed. A crucial requirement for the reduced model is that it be able to replicate the output quantities of interest (i.e. the observables) of the PDE simulation over a wide range of initial conditions, so that it may serve as a surrogate of the high fidelity PDE simulation during inversion.

One popular method for generating a reduced model is through a projection basis (for example, by proper orthogonal decomposition in conjunction with the method of snapshots). To build such a reduced order model, one typically constructs a training set by sampling the space of (discretized) initial conditions. When this space is high-dimensional, the problem of adequately sampling it quickly becomes intractable. Fortunately, for many *ill-posed* inverse problems, many components of the initial condition space have minimal or no effect on the output observables. This is particularly true when the observations are sparse. In this case, it is likely that an effective reduced model can be generated with few sample points. The question is how to locate these sample points.

Here, we consider the case of a linear forward problem, and propose that the sample points be associated with dominant eigenvectors of the Hessian matrix of the misfit function. This matrix maps inputs (initial conditions) to outputs (observables), and its dominant eigenvectors represent initial condition components that are most identifiable from observable data. Thus, one expects these eigenvectors to serve as good sample points for constructing the reduced model.

In Section 2, we describe the model reduction framework we consider, and in Section 3 justify the choice of the dominant eigenvectors of the Hessian by relating it to solution of a certain greedy optimization problem to locate the best sample points. Section 4 illustrates the methodology via application to a data assimilation inverse problem involving transport of an atmospheric contaminant.

2 Reduced-Order Dynamical Systems

Consider the general linear initial-value problem

$$x(k+1) = Ax(k), \quad k = 0, 1, \dots, T-1, \quad (1)$$

$$y(k) = Cx(k), \quad k = 0, 1, \dots, T, \quad (2)$$

$$x(0) = x_0, \quad (3)$$

where $x(k) \in \mathbb{R}^N$ is the system state at time t_k , the vector x_0 contains the specified initial state, and we consider a time horizon from $t = 0$ to $t = t_T$. The vector $y(k) \in \mathbb{R}^Q$ contains the Q system outputs at time t_k . In general, we are interested in systems of the form (1)–(3) that result from spatial and temporal discretization of PDEs. In this case, the dimension of the system, N , is very large and the matrices $A \in \mathbb{R}^{N \times N}$ and $C \in \mathbb{R}^{Q \times N}$ result from the chosen spatial and temporal discretization methods.

A reduced-order model of (1)–(3) can be derived by assuming that the state $x(k)$ is represented as a linear combination of n basis vectors,

$$\hat{x}(k) = Vx_r(k), \tag{4}$$

where $\hat{x}(k)$ is the reduced model approximation of the state $x(k)$ and $n \ll N$. The projection matrix $V \in \mathbb{R}^{N \times n}$ contains as columns the orthonormal basis vectors V_i , i.e., $V = [V_1 \ V_2 \ \dots \ V_n]$, and the reduced-order state $x_r(k) \in \mathbb{R}^n$ contains the corresponding modal amplitudes for time t_k . Using the representation (4) together with a Galerkin projection of the discrete-time system (1)–(3) onto the space spanned by the basis V yields the reduced-order model with state x_r and output y_r ,

$$x_r(k + 1) = A_r x_r(k), \quad k = 0, 1, \dots, T - 1, \tag{5}$$

$$y_r(k) = C_r x_r(k), \quad k = 0, 1, \dots, T, \tag{6}$$

$$x_r(0) = V^T x_0, \tag{7}$$

where $A_r = V^T A V$ and $C_r = C V$.

For convenience of notation, we write the discrete-time system (1)–(3) in matrix form as

$$\mathbf{A}\mathbf{x} = \mathbf{F}\mathbf{x}_0, \quad \mathbf{y} = \mathbf{C}\mathbf{x}, \tag{8}$$

where $\mathbf{x} = [x(0)^T \ x(1)^T \ \dots \ x(T)^T]^T$, $\mathbf{y} = [y(0)^T \ y(1)^T \ \dots \ y(T)^T]^T$, and the matrices \mathbf{A} , \mathbf{F} , and \mathbf{C} are appropriately defined functions of A and C . Similarly, the reduced-order model (5)–(7) can be written in matrix form as

$$\mathbf{A}_r \mathbf{x}_r = \mathbf{F}_r \mathbf{x}_0, \quad \mathbf{y}_r = \mathbf{C}_r \mathbf{x}_r, \tag{9}$$

where \mathbf{x}_r , \mathbf{y}_r , \mathbf{A}_r , and \mathbf{C}_r are defined analogously to \mathbf{x} , \mathbf{y} , \mathbf{A} , and \mathbf{C} but with the appropriate reduced-order quantities, and $\mathbf{F}_r = [V \ 0 \ \dots \ 0]^T$.

In many cases, we are interested in rapid identification of initial conditions from sparse measurements of the states over a time horizon; we thus require a reduced-order model that will provide accurate outputs for any initial condition contained in some set \mathcal{X}_0 . Using the projection framework described above, the task therefore becomes one of choosing an appropriate basis V so that the error between full-order output \mathbf{y} and the reduced-order output \mathbf{y}_r is small for all initial conditions of interest.

3 Hessian-Based Model Reduction

To determine the reduced model, we must identify a set of initial conditions to be sampled. At each selected initial condition, a forward simulation is performed to generate a set of states, commonly referred to as snapshots, from which the reduced basis is formed. The key question is then how to identify important initial conditions that should be sampled. Our approach is motivated by the greedy algorithm of [5], which proposed an adaptive approach to determine the parameter locations at which samples are drawn to form a reduced basis. The greedy algorithm adaptively selects these snapshots by finding the location in parameter–time space where the error between the full-order and reduced-order models is maximal, updating the basis with information gathered from this sample location, forming a new reduced model, and repeating the process.

In the case of the initial-condition problem, the greedy approach amounts to sampling at the initial condition $x_0^* \in \mathcal{X}_0$ that *maximizes* the error between the full and reduced-order outputs. For this formulation, the only restriction that we place on the set \mathcal{X}_0 is that it contain vectors of unit length. This prevents unboundedness in the optimization problem, since otherwise the error in the reduced system could be made arbitrarily large.

The key step in the greedy sampling approach is thus finding the worst-case initial condition x_0^* , which can be achieved by solving the optimization problem,

$$x_0^* = \arg \max_{x_0 \in \mathcal{X}_0} (\mathbf{y} - \mathbf{y}_r)^T (\mathbf{y} - \mathbf{y}_r) \tag{10}$$

$$\text{where } \mathbf{A}\mathbf{x} = \mathbf{F}x_0, \tag{11}$$

$$\mathbf{y} = \mathbf{C}\mathbf{x}, \tag{12}$$

$$\mathbf{A}_r\mathbf{x}_r = \mathbf{F}_rx_0, \tag{13}$$

$$\mathbf{y}_r = \mathbf{C}_r\mathbf{x}_r. \tag{14}$$

Equations (10)-(14) define a large-scale optimization problem, which includes the full-scale dynamics as constraints. The linearity of the state equations can be exploited to eliminate the full-order and reduced-order states and yield an equivalent unconstrained optimization problem,

$$x_0^* = \arg \max_{x_0 \in \mathcal{X}_0} x_0^T H^e x_0, \tag{15}$$

where

$$H^e = (\mathbf{C}\mathbf{A}^{-1}\mathbf{F} - \mathbf{C}_r\mathbf{A}_r^{-1}\mathbf{F}_r)^T (\mathbf{C}\mathbf{A}^{-1}\mathbf{F} - \mathbf{C}_r\mathbf{A}_r^{-1}\mathbf{F}_r). \tag{16}$$

It can be seen that (15) is a quadratic unconstrained optimization problem with Hessian matrix $H^e \in \mathbb{R}^{N \times N}$. From (16), it can be seen that H^e is a symmetric positive semidefinite matrix. Since we are considering initial conditions of unit norm, the solution x_0^* maximizes the Rayleigh quotient; therefore, the solution of (15) is given by the eigenvector corresponding to the largest eigenvalue of H^e .

This eigenvector is the initial condition for which the error in reduced model output prediction is largest.

Rather than constructing a reduced model at every greedy iteration, and determining the dominant eigenvector of the resulting error Hessian H_e , an efficient one-shot algorithm can be constructed by computing the dominant eigenmodes of the Hessian matrix

$$H = (\mathbf{CA}^{-1}\mathbf{F})^T (\mathbf{CA}^{-1}\mathbf{F}). \quad (17)$$

Here, $H \in \mathbb{R}^{N \times N}$ is the Hessian matrix of the full-scale system, and does not depend on the reduced-order model. As before, H is a symmetric positive semi-definite matrix. It can be shown that, under certain assumptions, the eigenvectors of H with largest eigenvalues approximately solve the sequence of problems defined by (10)–(14) [3].

These ideas motivate the following basis-construction algorithm for the initial condition problem. We use the dominant eigenvectors of the Hessian matrix H to identify the initial-condition vectors that have the most significant contributions to the outputs of interest. These vectors are in turn used to initialize the full-scale discrete-time system to generate a set of state snapshots that are used to form the reduced basis (using, for example, the proper orthogonal decomposition).

4 Application: Model Reduction for 3D Contaminant Transport in an Urban Canyon

We demonstrate our model reduction method by applying it to a 3D airborne contaminant transport problem for which a solution is needed in real time. Intentional or unintentional chemical, biological, and radiological (CBR) contamination events are important national security concerns. In particular, if contamination occurs in or near a populated area, predictive tools are needed to rapidly and accurately forecast the contaminant spread to provide decision support for emergency response efforts. Urban areas are geometrically complex and require detailed spatial discretization to resolve the relevant flow and transport, making prediction in real-time difficult. Reduced-order models can play an important role in facilitating real-time turn-around, in particular on laptops in the field. However, it is essential that these reduced models be faithful over a wide range of initial conditions, since in principle any initial condition can be realized. Once a suitable reduced-order model has been generated, it can serve as a surrogate for the full model within an inversion framework to identify the initial conditions given sensor data (the full-scale case is discussed in [1]).

To illustrate the generation of a reduced-order model that is accurate for arbitrary high-dimensional initial conditions, we consider a three-dimensional urban canyon geometry occupying a (dimensionless) $15 \times 15 \times 15$ domain. Figure 1 shows the domain and buildings, along with locations of six sensors, all

placed at a height of 1.5. Contaminant transport is modeled by the advection-dispersion equation,

$$\frac{\partial w}{\partial t} + \mathbf{v} \cdot \nabla w - \kappa \nabla^2 w = 0 \quad \text{in } \Omega \times (0, t_f), \tag{18}$$

$$w = 0 \quad \text{on } \Gamma_D \times (0, t_f), \tag{19}$$

$$\frac{\partial w}{\partial n} = 0 \quad \text{on } \Gamma_N \times (0, t_f), \tag{20}$$

$$w = w_0 \quad \text{in } \Omega \text{ for } t = 0, \tag{21}$$

where w is the contaminant concentration, \mathbf{v} is the velocity vector field, κ is the diffusivity, t_f is the time horizon of interest, and w_0 is the given initial condition. Γ_D and Γ_N are respectively the portions of the domain boundary over which Dirichlet and Neumann boundary conditions are applied. Eq. (18) is discretized in space using an SUPG finite element method with linear tetrahedra, while the implicit Crank-Nicolson method is used to discretize in time. Homogeneous Dirichlet boundary conditions are specified for the concentration on the inflow boundary, $\bar{x} = 0$, and the ground, $\bar{z} = 0$. Homogeneous Neumann boundary conditions are specified for the concentration on all other boundaries.

The velocity field, \mathbf{v} , required in (18) is computed by solving the steady laminar incompressible Navier-Stokes equations, also discretized with SUPG-stabilized linear tetrahedra. No-slip conditions, i.e. $\mathbf{v} = 0$, are imposed on the building faces and the ground $\bar{z} = 0$. The velocity at the inflow boundary $\bar{x} = 0$ is taken as known and specified in the normal direction as

$$v_x(z) = v_{\max} \left(\frac{z}{z_{\max}} \right)^{0.5},$$

with $v_{\max} = 3.0$ and $z_{\max} = 15$, and zero tangentially. On the outflow boundary $\bar{x} = 15$, a traction-free (Neumann) condition is applied. On all other boundaries ($\bar{y} = 0, \bar{y} = 15, \bar{z} = 15$), we impose a combination of no flow normal to the boundary and traction-free tangent to the boundary. The spatial mesh for the full-scale system contains 68,921 nodes and 64,000 tetrahedral elements. For both basis creation and testing, a final non-dimensional time $t_f = 20.0$ is used, and discretized over 200 timesteps. The Peclet number based on the maximum inflow velocity and domain dimension is $Pe=900$. The PETSc library [2] is used for all implementation.

Figure 2 illustrates a sample forward solution. The test initial condition used in this simulation, meant to represent the system state just after a contaminant release event, was constructed using a Gaussian function with a peak magnitude of 100 centered at a height of 1.5. For comparison with the full system, a reduced model was constructed based on the dominant Hessian eigenvector algorithm discussed in the previous section, with $p = 31$ eigenvector initial conditions and $n = 137$ reduced basis vectors (these numbers were determined based on eigenvalue decay rates). Eigenvectors were computed using the Arnoldi eigensolver within the SLEPc package [4], which is built on PETSc. Figure 3 shows a comparison of the full and reduced time history of concentration at each output

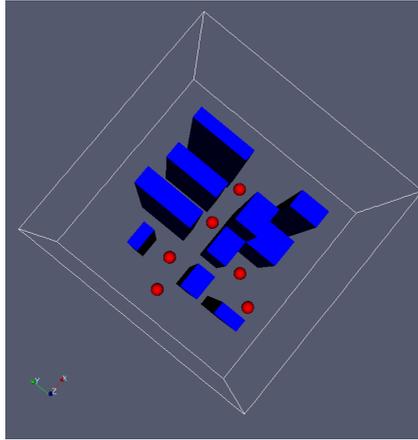


Fig. 1. Building geometry and locations of outputs for the 3-D urban canyon problem

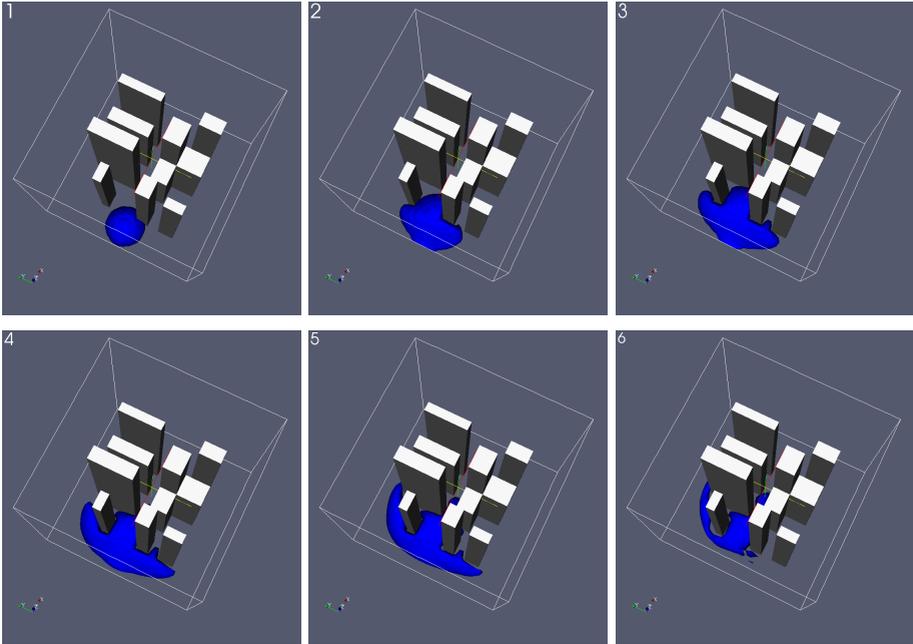


Fig. 2. Transport of contaminant concentration through urban canyon at six instants in time, beginning with the initial condition shown in upper left

location. There is no discernible difference between the two. The figure demonstrates that a reduced system of size $n = 137$, which is solved in a matter of seconds on a desktop, can accurately replicate the outputs of the full-scale system of size $N = 65,600$. We emphasize that the (offline) construction of the

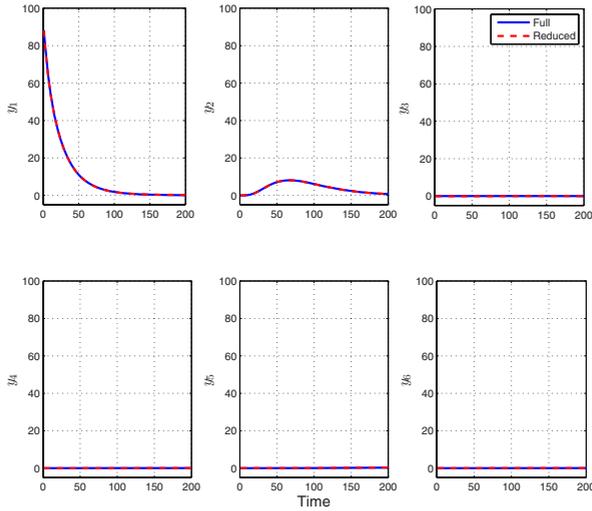


Fig. 3. Full (65,600 states) and reduced (137 states) model contaminant predictions at the six sensor locations for urban canyon example

reduced-order model targets only the specified outputs, and otherwise has no knowledge of the initial conditions used in the test of Figure 3.

References

1. V. Akçelik, G. Biros, A. Draganescu, O. Ghattas, J. Hill, and B. van Bloemen Waanders. Dynamic data-driven inversion for terascale simulations: Real-time identification of airborne contaminants. In *Proceedings of SC2005, Seattle, WA, 2005*.
2. S. Balay, K. Buschelman, V. Eijkhout, W. Gropp, D. Kaushik, M. Knepley, L. McInnes, B. Smith, and H. Zhang. PETSc users manual. Technical Report ANL-95/11 - Revision 2.1.5, Argonne National Laboratory, 2004.
3. O. Bashir. Hessian-based model reduction with applications to initial condition inverse problems. Master's thesis, MIT, 2007.
4. V. Hernandez, J. Roman, and V. Vidal. SLEPc: A scalable and flexible toolkit for the solution of eigenvalue problems. *ACM Transactions on Mathematical Software*, 31(3):351–362, sep 2005.
5. K. Veroy, C. Prud'homme, D. Rovas, and A. Patera. *A posteriori* error bounds for reduced-basis approximation of parametrized noncoercive and nonlinear elliptic partial differential equations. AIAA Paper 2003-3847, Proceedings of the 16th AIAA Computational Fluid Dynamics Conference, Orlando, FL, 2003.