



### Science Arts & Métiers (SAM)

is an open access repository that collects the work of Arts et Métiers Institute of Technology researchers and makes it freely available over the web where possible.

This is an author-deposited version published in: <https://sam.ensam.eu>  
Handle ID: <http://hdl.handle.net/10985/17796>

#### To cite this version :

Steven L. BRUNTON, Jean-Christophe LOISEAU - Constrained sparse Galerkin regression -  
Journal of Fluid Mechanics - Vol. 838, p.42-67 - 2018

Any correspondence concerning this service should be sent to the repository

Administrator : [scienceouverte@ensam.eu](mailto:scienceouverte@ensam.eu)



# Constrained Sparse Galerkin Regression

J.-Ch. Loiseau<sup>1</sup> and S. L. Brunton<sup>2</sup>

<sup>1</sup>Laboratoire DynFluid, Arts et Métiers ParisTech, 75013 Paris, France

<sup>2</sup>Department of Mechanical Engineering, University of Washington, Seattle, WA 98195, USA

## Abstract

Although major advances have been achieved over the past decades for the reduction and identification of linear systems, deriving nonlinear low-order models still is a challenging task. In this work, we develop a new data-driven framework to identify nonlinear reduced-order models of a fluid by combining dimensionality reductions techniques (*e.g.* proper orthogonal decomposition) and sparse regression techniques from machine learning. In particular, we extend the sparse identification of nonlinear dynamics (SINDy) algorithm to enforce physical constraints in the regression, namely energy-preserving quadratic nonlinearities. The resulting models, hereafter referred to as *Galerkin regression* models, incorporate many beneficial aspects of Galerkin projection, but without the need for a full-order or high-fidelity solver to project the Navier-Stokes equations. Instead, the most parsimonious nonlinear model is determined that is consistent with observed measurement data and satisfies necessary constraints. Galerkin regression models also readily generalize to include higher-order nonlinear terms that model the effect of truncated modes. The effectiveness of Galerkin regression is demonstrated on two different flow configurations: the two-dimensional flow past a circular cylinder and the shear-driven cavity flow. For both cases, the accuracy of the identified models compare favorably against reduced-order models obtained from a standard Galerkin projection procedure. Present results highlight the importance of cubic nonlinearities in the construction of accurate nonlinear low-dimensional approximations of the flow systems, something which cannot be readily obtained using a standard Galerkin projection of the Navier-Stokes equations. Finally, the entire code base for our constrained sparse Galerkin regression algorithm is freely available online.

## 1 Introduction

Fluid flows are characterised by high-dimensional, nonlinear dynamics that give rise to rich structures. Despite this apparent complexity, the dynamics often evolve on a low-dimensional attractor defined by a few dominant coherent structures that contain significant energy or are useful for control (Holmes et al., 2012). Given this property, one might then aim to derive or identify reduced-order models that reproduce qualitatively and quantitatively the dynamics of the full system. Over the past decades, identifying robust, accurate and efficient reduced-order models has thus become a central challenge in fluid dynamics and closed-loop flow control (Brunton and Noack, 2015; Fabbiane et al., 2014; Rowley and Dawson, 2016; Sipp and Schmid, 2016).

Many traditional model reduction techniques are analytical. They rely on prior knowledge of the Navier-Stokes equations and the existence of a high-fidelity solver to project onto an orthogonal

basis of modes, resulting in a dynamical system in terms of the coefficients of this expansion basis. These modes may come from a classical expansion, such as Fourier modes, or they may be data-driven, as in the proper orthogonal decomposition (POD) (Berkooz et al., 1993; Sirovich, 1987). In the latter case, the model-reduction may be considered a hybrid approach, mixing knowledge of the physics with empirical modes obtained from measurement data. Control-theoretic extensions, such as balanced POD (BPOD) (Rowley, 2005; Willcox and Peraire, 2002), have also been widely applied for closed-loop flow control (Bagheri et al., 2009; Ilak and Rowley, 2008; Illingworth et al., 2010). Although such approaches to model reduction have been widely successful for linear systems, as described in the recent review by Rowley and Dawson (2016) and references therein, they have been applied with only limited success to obtain low-order approximations of nonlinear systems, mostly on flow oscillators. One can cite for instance the seminal work of Noack et al. (2003) and Tadmor et al. (2010) wherein the authors have shown that such reduced-order models obtained from a Galerkin projection can reproduce the transients and non-linear dynamics of the von Kàrmàn vortex shedding past a two-dimensional cylinder provided the projection basis includes a *shift mode* quantifying the distortion between the linearly unstable base flow and marginally stable mean flow. Recently, Semaan et al. (2016) have extended the reduced-order modeling strategy of Noack et al. (2003) to include the effect of control actuation for the flow around a high-lift configuration airfoil.

In contrast, data-driven approaches are becoming increasingly popular and encompass a large variety of different techniques such as the eigensystem realisation algorithm (ERA) (Juang and Pappa, 1985), dynamic mode decomposition (DMD) (Kutz et al., 2016; Rowley et al., 2009; Schmid, 2010), Koopman theory (Mezić, 2005, 2013) and variants (Tu et al., 2014; Williams et al., 2015), cluster reduced order modeling (CROM) (Kaiser et al., 2014), and network analysis of fluids (Nair and Taira, 2015). Recent advances in machine learning are also greatly expanding the ability to extract governing dynamics purely from data. In particular advanced regression methods from statistics, such as genetic programming or sparse regression, are driving new algorithms that identify nonlinear dynamics from measurements of complex systems. Bongard and Lipson (2007) and Schmidt and Lipson (2009) introduced nonlinear system identification based on genetic programming, which has been used in numerous practical applications in aerospace engineering, the petroleum industry, and in finance. More recently, Brunton et al. (2016b) have proposed a system identification approach based on sparse regression known as the sparse identification of nonlinear dynamics (SINDy). Following the principle of Ockham’s razor, the SINDy algorithm rests on the assumption that there are only a few important terms that govern the dynamics of a given system, so that the equations are sparse in the space of possible functions. Sparse regression is then used to determine the fewest terms in a dynamical system required to accurately represent the data. The resulting models are parsimonious, balancing model complexity with descriptive power while avoiding overfitting.

Most of these regression techniques can be recast into a convex minimisation problem and their solution can be obtained using a number of efficient algorithms available in different libraries such as CVXOPT (Andersen et al., 2013). However, a major drawback of regression-based methods is the possible loss of existing symmetries in the governing equations which may otherwise be included in the physics-based Galerkin projection methods described previously (Balajewicz et al., 2013; Carlberg et al., 2015). A notable exception is the physics-constrained multi-level quadratic regression used to identify models in climate and turbulence (Majda and Harlim, 2012). Starting from the original SINDy algorithm (Brunton et al., 2016b), a system identification technique based

on sparse regression, we propose in this work a new implementation of the algorithm which allows the user to include physical constraints such as energy-preserving nonlinearities or to enforce symmetries in the identified equations. The resulting algorithm relies on the use of constrained least squares (Golub and Van Loan, 2012) to incorporate additional constraints in the SINDy algorithm for the sparse identification of the underlying low-dimensional dynamical system. The ability of the present system identification technique, hereafter named *sparse Galerkin regression*, is demonstrated on two different flow configurations, namely the emblematic two-dimensional cylinder flow and the shear-driven cavity flow. The manuscript is organised as follows: §2.1 provides the reader with a quick introduction to the original SINDy algorithm, while the new algorithm is presented in §2.2 and illustrated on a toy model in §2.3. The physical constraints used in this work are discussed in §3, while the two flow configurations considered herein are presented in §4. The different low-dimensional systems identified are compared against standard Galerkin projection in §5. Finally, §6 summarises our key findings and provide the reader with possible extensions to this work.

## 2 Constrained sparse identification

Here we discuss the core mathematical and algorithmic framework used to identify nonlinear reduced-order models from data. The proposed Galerkin regression method is based on a heavily modified version of the sparse identification of nonlinear dynamics (SINDy) method (Brunton et al., 2016b). The original SINDy algorithm is introduced in § 2.1, and the new modifications to include physical constraints, such as energy conservation, known eigenvalues, or symmetries, are discussed in § 2.2. Implementation details for both algorithms are presented to promote reproducibility; in addition, code is freely available online (<https://github.com/loiseaujc/SINDy>). Finally, the inclusion of constraints is demonstrated on the Lorenz system as an illustrative example in § 2.3. Specific constraints that are used to enforce energy conservation are derived later in § 3.

### 2.1 Sparse identification of nonlinear dynamics (SINDy)

Identifying dynamical systems models from data has been a central challenge in mathematical physics, with a particularly rich history in fluid dynamics. Typically, the form of the dynamical systems model identified is either constrained via prior knowledge of the governing equations, as in Galerkin projection, or a small handful of heuristic models are posited and parameters are optimized to match the data. Simultaneously identifying the structure and parameters of a model from data is considerably more challenging, as there are combinatorially many possible model structures.

The sparse identification of nonlinear dynamics (SINDy) algorithm (Brunton et al., 2016b) bypasses the intractable brute force search through all possible model structures, leveraging the observation that many dynamical systems

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) \tag{1}$$

have dynamics  $\mathbf{f}$  that are sparse in the space of possible right-hand side functions. It is then possible to solve for the relevant terms that are active in the dynamics using a convex  $\ell_1$ -regularized regression that penalizes the number of terms in the dynamics and scales well to large problems.

First, time-series data is collected from Eq. (1) and formed into a data matrix:

$$\mathbf{X} = [\mathbf{x}(t_1) \quad \mathbf{x}(t_2) \quad \cdots \quad \mathbf{x}(t_m)]^T \tag{2}$$

where  $T$  denotes the matrix transpose. A similar matrix of derivatives is formed:

$$\dot{\mathbf{X}} = [\dot{\mathbf{x}}(t_1) \quad \dot{\mathbf{x}}(t_2) \quad \cdots \quad \dot{\mathbf{x}}(t_m)]^T. \quad (3)$$

In practice, this may be computed directly from the data in  $\mathbf{X}$ ; for noisy data, the total-variation regularized derivative tends to provide numerically robust derivatives (Chartrand, 2011).

Based on the data in  $\mathbf{X}$ , a library of candidate nonlinear functions  $\Theta(\mathbf{X})$  is constructed:

$$\Theta(\mathbf{X}) = [\mathbf{1} \quad \mathbf{X} \quad \mathbf{X}^2 \quad \cdots \quad \mathbf{X}^d \quad \cdots \quad \sin(\mathbf{X}) \quad \cdots]. \quad (4)$$

Here, the matrix  $\mathbf{X}^d$  denotes a matrix with column vectors given by all possible time-series of  $d$ -th degree polynomials in the state  $\mathbf{x}$ .

The dynamical system in Eq. (1) may now be represented in terms of the data matrices in Eqs. (3) and (4) as

$$\dot{\mathbf{X}} = \Theta(\mathbf{X})\Xi. \quad (5)$$

Each column  $\Xi_k$  in  $\Xi$  is a vector of coefficients determining the active terms in the  $k$ -th row equation in Eq. (1). A parsimonious model will provide an accurate model fit in Eq. (5) with as few terms as possible in  $\Xi$ . Such a model may be identified using a convex  $\ell_1$ -regularized sparse regression:

$$\Xi_k = \operatorname{argmin}_{\Xi'_k} \|\dot{\mathbf{X}}_k - \Theta(\mathbf{X})\Xi'_k\|_2 + \lambda \|\Xi'_k\|_1. \quad (6)$$

Here,  $\dot{\mathbf{X}}_k$  is the  $k$ -th column of  $\dot{\mathbf{X}}$ . Sparse regression, such as the LASSO (Tibshirani, 1996) or the sequential thresholded least-squares algorithm used in SINDy, improves the numerical robustness of this identification for noisy overdetermined problems, in contrast to earlier methods (Wang et al., 2011) that used compressed sensing (Candès, 2006; Donoho, 2006).

The sparse vectors  $\Xi_k$  may be synthesized into a nonlinear dynamical system model:

$$\dot{x}_k = \Theta(\mathbf{x})\Xi_k. \quad (7)$$

Note that  $x_k$  is the  $k$ -th element of  $\mathbf{x}$  and  $\Theta(\mathbf{x})$  is a row vector of symbolic functions of  $\mathbf{x}$ , as opposed to the data matrix  $\Theta(\mathbf{X})$ .

Identifying the most parsimonious nonlinear model by applying sparse regression in the library  $\Theta$  is a convex procedure. The alternative approach, which involves regression onto every possible sparse nonlinear structure, constitutes an intractable brute-force procedure. SINDy bypasses this combinatorial search with modern convex optimization and machine learning. It is interesting to note that if  $\Theta(\mathbf{X})$  consists only of linear terms, and if we remove the sparsity promoting term by setting  $\lambda = 0$ , then this algorithm reduces to the dynamic mode decomposition Kutz et al. (2016); Rowley et al. (2009); Schmid (2010).

Recent extension to SINDy enable the identification of nonlinear differential equations with rational function nonlinearities by reformulating the problem as an implicit differential equation and solving for the active terms by finding the sparsest vector in the null space of an augmented library containing functions of the state and derivative terms (Mangan et al., 2016). SINDy has also been generalized to identify partial differential equations from data (Rudy et al., 2016), and has been extended to include inputs and control (Brunton et al., 2016a).

## 2.2 Constrained sparse identification

It has been shown in §2.1 that, within the SINDy framework, the identification problem can be cast as a convex optimisation problem where the sparsity of the solution  $\Xi$  can be promoted using an  $l_1$  regularized regression. Alternatively, sparsity can also be promoted by using the sequential thresholded least-squares algorithm as in Brunton et al. (2016b). In this case, the convex minimisation problem can be re-written as

$$\begin{aligned} \min_{\Xi} \|\Theta(\mathbf{X})\Xi - \dot{\mathbf{X}}\|_2^2 \\ \text{subject to } \mathbf{C}\xi = \mathbf{d} \end{aligned} \quad (8)$$

where  $\xi = \Xi(\cdot)$  is the vectorized form of the sparse matrix of coefficients, and where  $\mathbf{C}\xi = \mathbf{d}$  are linear equality constraints, which can be used to enforce that some entries of  $\xi$  are equal to zero. The minimisation problem is then solved iteratively. After an initial least-squares regression, the thresholding is performed as follows: if  $|\xi_i|$  is smaller than  $\lambda$  (the sparsity knob) times the mean of the absolute value of the non-zero entries of  $\xi$ , then an additional row is added to the constraint matrix  $\mathbf{C}$  to enforce  $\xi_i = 0$ . Two or three iterations of this small variation of the sequential thresholded least-squares algorithm are usually sufficient to ensure convergence of the constrained minimization procedure. The sparsity parameter  $\lambda$  should be chosen to promote parsimonious models that strike a balance between accuracy and complexity to avoid overfitting the data. More details on this choice are presented in Appendix B.

From a practical point of view, each iteration of (8) can be recast as an unconstrained problem by using an augmented functional formulation where the constraints are imposed by means of Lagrange multipliers. The resulting unconstrained minimisation problem then reads

$$\min_{\xi, \mathbf{z}} \|\Theta(\mathbf{X})\Xi - \dot{\mathbf{X}}\|_2^2 + \mathbf{z}^T(\mathbf{C}\xi - \mathbf{d}). \quad (9)$$

Given the choice of our augmented functional, it can easily be shown that the optimal solution  $\xi$  that satisfies the constraints is also solution to the following Karush-Kuhn-Tucker (KKT) equations

$$\begin{bmatrix} 2\hat{\Theta}(\mathbf{X})^T \hat{\Theta}(\mathbf{X}) & \mathbf{C}^T \\ \mathbf{C} & 0 \end{bmatrix} \begin{bmatrix} \xi \\ \mathbf{z} \end{bmatrix} = \begin{bmatrix} 2\hat{\Theta}(\mathbf{X})^T \dot{\mathbf{X}}(\cdot) \\ \mathbf{d} \end{bmatrix}, \quad (10)$$

where  $\hat{\Theta}(\mathbf{X})$  is a diagonal matrix consisting of  $n$  copies of  $\Theta(\mathbf{X})$ ,  $\mathbf{X}(\cdot)$  is the vectorized form of  $\mathbf{X}$  (same as the vectorization of  $\Xi$  into  $\xi = \Xi(\cdot)$ ), and  $n$  is the dimension of  $\mathbf{x}$ . This matrix equation for constrained least-squares is the counterpart to the ordinary least-squares normal equations. It has a unique solution if  $\mathbf{C}$  has full row-rank and  $[\hat{\Theta}(\mathbf{X}) \quad \mathbf{C}]^T$  has full column-rank.

Interestingly, the linear equality constraints  $\mathbf{C}\xi = \mathbf{d}$  do not have to be used for the sole purpose of sparsity promotion. Indeed, these can also be used to enforce additional user-provided constraints such as an *a priori* known value of a given entry  $\xi_i$  or to impose some linear relationship between the entries of  $\xi$  to mimic a given physical process, see §2.3 for a simple illustration. Specific constraints required to conserve energy in a fluid are derived later in § 3.

### 2.3 Illustration of constrained sparse identification on the Lorenz system

Following Brunton et al. (2016b), let us first illustrate how to formulate user-provided constraints using the Lorenz system (Lorenz, 1963). This dynamical system, derived by Edward Lorenz in 1963, is notable for having chaotic solutions for certain parameter values and initial conditions. It reads

$$\begin{aligned}\dot{x} &= \sigma(y - x) \\ \dot{y} &= x(\rho - z) - y \\ \dot{z} &= xy - \beta z.\end{aligned}\tag{11}$$

Figure 1 depicts the evolution in time of  $\mathbf{x}(t) = [x(t), y(t), z(t)]^T$  for a given set of parameters  $\sigma$ ,  $\rho$  and  $\beta$ . These signals, as well as their derivatives (not shown), will serve as the input data for the constrained system identification. For that purpose, the library  $\Theta(\mathbf{x})$  used in the identification process is defined as  $P_2(\mathbf{x})$ , *i.e.* all the polynomials of degree 2 or less in the entries of  $\mathbf{x}$  such that

$$\Theta(\mathbf{x}) = [1 \ x \ y \ z \ x^2 \ xy \ xz \ y^2 \ yz \ z^2].\tag{12}$$

Up to 30 different coefficients thus need to be identified, 10 per equation. Let us assume furthermore that, in the  $x$ -equation, we know beforehand that  $\sigma = 10$ . The constrained optimisation problem on which SINDY relies then reads

$$\begin{aligned}\min_{\xi} \quad & \|\dot{\mathbf{X}} - \Theta(\mathbf{X})\xi\|_2^2 \\ \text{subject to} \quad & \xi_3 = 10 \\ & \xi_2 + \xi_3 = 0.\end{aligned}\tag{13}$$

From a practical point of view, these equality constraints are passed to CVXOPT as  $\mathbf{C}\xi = \mathbf{d}$ , where  $\mathbf{C}$  is a  $2 \times 30$  matrix and  $\mathbf{d}$  a vector given by

$$\mathbf{C} = \begin{bmatrix} 0 & 0 & 1 & 0 & \cdots & 0 \\ 0 & 1 & 1 & 0 & \cdots & 0 \end{bmatrix} \text{ and } \mathbf{d} = [10 \ 0]^T.\tag{14}$$

Using a suitable sparsity knob, the system identified by the constrained SINDy algorithm finally reads

$$\begin{aligned}\dot{x} &= 10(y - x) \\ \dot{y} &= x(27.99 - 0.999z) - 0.998y \\ \dot{z} &= 0.999xy - 2.666z.\end{aligned}\tag{15}$$

The coefficients of the identified system are close to the original ones, which were set to  $\sigma = 10$ ,  $\rho = 28$  and  $\beta = 8/3$ . The time-evolution given by this identified system is depicted in figure 1 along with the original signals and those given by a system identified using the original (unconstrained) SINDy algorithm. It can be seen that the trajectory of the system identified using constrained SINDy remains closer to that of the original system compared to the trajectory predicted by the system identified using the original SINDy. The effects of adding constraints are even more pronounced in fluid systems where energy conservation may be enforced if certain constraints on the quadratic nonlinearities are satisfied, as discussed in § 3 and § 5.

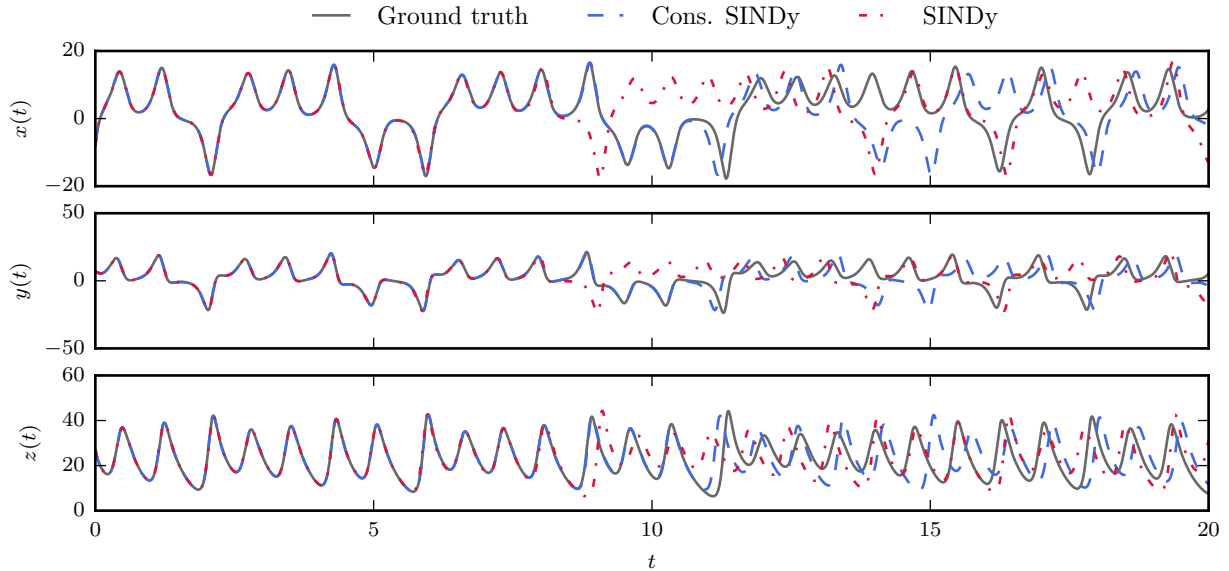


Figure 1: Dataset (grey) used for the constrained sparse identification of the Lorenz system and prediction obtained using the model identified by the constrained SINDy algorithm (blue dashed) and the original SINDy algorithm (red dashed).

### 3 Deriving the constraints

The Navier-Stokes equations governing the dynamics of the perturbation  $\mathbf{u}$  evolving on top of the base flow  $\mathbf{U}_b$  are given by

$$\frac{\partial \mathbf{u}}{\partial t} = -(\mathbf{U}_b \cdot \nabla) \mathbf{u} - (\mathbf{u} \cdot \nabla) \mathbf{U}_b - \nabla p + \frac{1}{Re} \nabla^2 \mathbf{u} - (\mathbf{u} \cdot \nabla) \mathbf{u} \quad (16)$$

$$\nabla \cdot \mathbf{u} = 0 \quad (17)$$

where  $\mathbf{U}_b$  is the base flow velocity field,  $\mathbf{u}$  is the perturbation velocity field and  $p$  the corresponding pressure. The aim of reduced-order modeling is to construct/derive/identify a low-dimensional system of the form

$$\frac{d\mathbf{a}}{dt} = \tilde{\mathcal{L}}\mathbf{a} + \tilde{\mathcal{N}}(\mathbf{a})\mathbf{a} \quad (18)$$

where  $\tilde{\mathcal{L}}$  and  $\tilde{\mathcal{N}}(\mathbf{a})\mathbf{a}$  are low-dimensional approximation of the linearised Navier-Stokes operator and of the quadratic nonlinear term, respectively, and where the entries of  $\mathbf{a}$  are the degrees of freedom of the reduced-order model. For the reduced-order model (18) to be a good approximation of its high-dimensional counterpart, the former needs to have the same physical properties as the latter. While this is expected to be true when the reduced-order model is derived based on a Galerkin projection, these properties need to be actively enforced when a system identification approach as SINDy is used.



### 3.1 Constraining the quadratic nonlinear term

The nonlinear Navier-Stokes equations (17) are partial differential equations characterised by the quadratic nonlinear term  $-(\mathbf{u} \cdot \nabla)\mathbf{u}$ . It can be shown that

$$\int_{\Omega} \mathbf{u} \cdot (\mathbf{u} \cdot \nabla)\mathbf{u} \, d\Omega = 0 \quad (19)$$

where the boundary terms resulting from the integration by parts are assumed to be small enough and can thus be neglected for the sake of simplicity. The contribution of the quadratic nonlinear term to the total energy of the perturbation is zero: it is an energy-preserving nonlinearity, its role being only to redistribute the perturbation's energy along the different lengthscales of the problem.

Given that our projection basis contains the POD modes, their amplitudes  $a_i(t)$  are directly related to the kinetic energy of the perturbation. The constraint required in our system identification for the low-dimensional quadratic nonlinear term to be energy-preserving is thus

$$\mathbf{a} \cdot \tilde{\mathcal{N}}(\mathbf{a})\mathbf{a} = 0. \quad (20)$$

Expanding (20) in terms of the regression coefficients  $\xi$  yields

$$0 = \mathbf{a}^T \begin{bmatrix} \xi_4^{(a_1)} a_1^2 & \xi_5^{(a_1)} a_1 + \xi_7^{(a_1)} a_2 & \xi_6^{(a_1)} a_1 + \xi_9^{(a_1)} a_{\Delta} \\ \xi_4^{(a_2)} a_1 + \xi_5^{(a_2)} a_2 & \xi_7^{(a_2)} a_2 & \xi_8^{(a_2)} a_2 + \xi_9^{(a_2)} a_{\Delta} \\ \xi_4^{(a_{\Delta})} a_1 + \xi_6^{(a_{\Delta})} a_{\Delta} & \xi_7^{(a_{\Delta})} a_y + \xi_8^{(a_{\Delta})} a_{\Delta} & \xi_9^{(a_{\Delta})} a_{\Delta} \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_{\Delta} \end{bmatrix} + \mathbf{a}^T \begin{bmatrix} \xi_8^{(a_1)} a_2 a_{\Delta} \\ \xi_6^{(a_2)} a_1 a_{\Delta} \\ \xi_5^{(a_{\Delta})} a_1 a_2 \end{bmatrix} \quad (21)$$

For (21) to hold, the matrix involved in the first term is required to be skew-symmetric, while the second term implies  $\xi_8^{(a_1)} + \xi_6^{(a_2)} + \xi_5^{(a_{\Delta})} = 0$ . Overall, this gives rise to ten different linear equality constraints which induce a coupling of the different ordinary differential equations governing the evolution of  $a_1$ ,  $a_2$  and  $a_{\Delta}$ .

### 3.2 What about higher order nonlinearities?

Reduced-order modelling based on Galerkin projection usually requires a relatively large projection basis. Despite the very low effective dimensionality of the cylinder flow at  $Re = 100$ , Noack et al. (2003) demonstrated the need to include the first eight POD modes along with the shift mode for the reduced-order model to provide a relatively faithful approximation of the original high-dimensional dynamics. Including the higher harmonic POD modes was deemed necessary in order to limit the energy overshoot otherwise observed during the nonlinear saturation process. Even though they might be required to prevent a non-physical behaviour of the reduced-order model, these higher harmonic modes have very low energy and limited dynamics of their own: they are essentially *enslaved* to the dominant POD modes. Using *adiabatic elimination* (Haken, 1983) or *center manifold reduction* (Carini et al., 2015; Wiggins, 2003), it is well known that these slaved modes can be reduced out of the problem, while their influence onto the driving modes can be accounted for by appropriately modifying the nonlinear terms, generally introducing higher-order nonlinearities. Such an approach to reduced-order modelling, which can be summarised as *derive-then-reduce*, can be used to reduce the eight-dimensional system derived by Noack et al. (2003) for the two-dimensional cylinder flow into one having only three degrees of freedom, *i.e.* the amplitude of the shift mode and that of the first two POD modes.

This derive-then-reduce approach is generally quite involved, requiring cumbersome calculations, particularly if the original Galerkin projection model has more than just a few degrees of freedom. However, in the present work, high-order nonlinearities modelling the influence of the truncated modes can be automatically incorporated in the identification process, with no additional post-analysis. For that purpose, the library  $\Theta(\mathbf{a})$  of admissible functions needs to be extended in order to include higher-order polynomials. Note, however, that it is unclear at the present time how to constrain these high-order nonlinearities to ensure that the identified model is physical, although the method is effective in practice without constraining the higher-order terms.

## 4 Flow configurations

To demonstrate the Galerkin regression framework, we consider two prototypical flow configurations, the incompressible flow past a circular cylinder and the shear-driven cavity flow. These flows have been selected because they are standard benchmark problems for modal analysis, model reduction, and control in the literature, and because they provide a balance between complexity and interpretability.

### 4.1 Cylinder flow

The first flow configuration considered in the present work is the two-dimensional incompressible viscous flow past a circular cylinder at  $Re = 100$ . This Reynolds number, based on the free-stream velocity  $U_\infty$ , the cylinder diameter  $D$  and the kinematic viscosity  $\nu$ , is well above the onset of vortex shedding (Schumm et al., 1994; Zebib, 1987) and below the onset of three-dimensional instabilities (Barkley and Henderson, 1996; Zhang et al., 1995). In the fluid dynamics community, a large body of literature exists in which this particular setup has been chosen to illustrate modal decomposition (Bagheri, 2013) and model identification techniques (Brunton et al., 2016b; Noack et al., 2003; Rowley and Dawson, 2016; Sengupta et al., 2015). This setup is thus a particularly compelling test case to illustrate our model identification strategy, as well as to draw connections and quantify its performance against other well-established techniques, namely Galerkin projection.

The dynamics of the flow are governed by the incompressible Navier-Stokes equations. These are solved numerically using the Nek 5000 spectral element solver (Fischer et al., 2008). The same computational domain as in Noack et al. (2003) has been considered. It extends from  $x_1 = -5$  up to  $x_1 = 15$  in the streamwise direction, and from  $x_2 = -5$  up to  $x_2 = 5$  in the spanwise direction. It is discretised using 1832 seventh-order spectral elements. The vorticity field of the linearly unstable fixed point  $\mathbf{U}_b$ , computed using the selective frequency damping approach (Åkervik et al., 2006), is shown in figure 2(b). Figure 2(a) and (c) also provide the eigenspectrum of the linearised Navier-Stokes operator and the vorticity field associated to the leading unstable eigenmode for the sake of completeness. Though this eigenmode is clearly related to vortex shedding, it is well known that both its spatial distribution and the frequency of the associated eigenvalue differ quite significantly from that of the non-linearly saturated von Kàrmàn vortex street (Barkley, 2006).

Given this linearly unstable base flow as initial condition, a direct numerical simulation has been run until a statistically steady-state has been achieved. The dynamics of the system on the final attractor are then equidistantly sampled using  $M = 1000$  velocity field snapshots with a sampling frequency about 30 times larger than the vortex shedding frequency (Noack et al., 2015). The shift

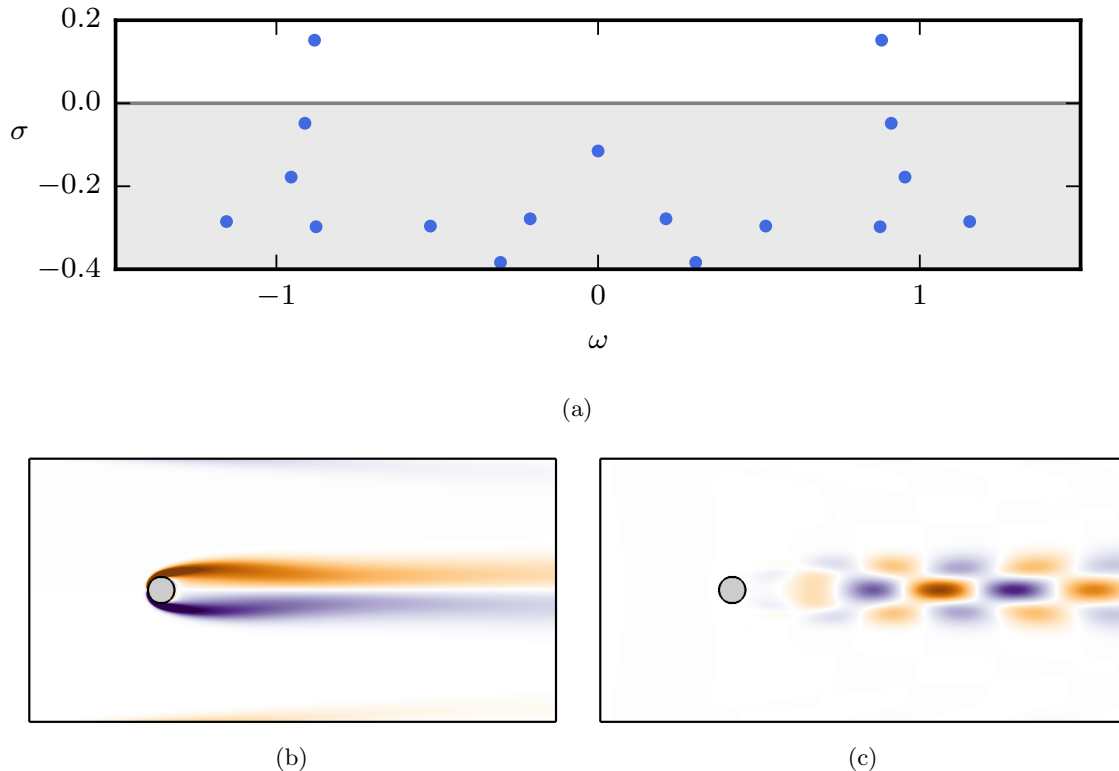


Figure 2: (a) Eigenspectrum of the linearised Navier-Stokes operator for the two-dimensional cylinder flow at  $Re = 100$ . Vorticity fields of (b) the base flow and (c) the leading linearly unstable eigenmode.

mode, denoted  $\mathbf{u}_\Delta$  and depicted in figure 3(a), quantifies the distortion between the unstable base flow equilibrium and the mean flow. It has been shown to be crucially important for POD-based reduced-order modeling (Noack et al., 2003; Tadmor et al., 2010). The snapshot POD method of Sirovich (1987) has then been used to extract the two most energetic modes  $\mathbf{u}_1$  and  $\mathbf{u}_2$ , depicted in figures 3(b) and 3(c), respectively. The evolution in time of the POD coefficients is shown in figure 4(a), while a projection of the system’s trajectory onto the  $a_1 - a_\Delta$  plane is depicted in figure 4(b), where  $a_1(t)$  is the amplitude of the POD mode  $\mathbf{u}_1$  and  $a_\Delta(t)$  the amplitude of the shift mode  $\mathbf{u}_\Delta$ . These signals and their time derivatives (not shown) form the training dataset used to identify the models in §5.1.

## 4.2 Shear-driven cavity flow

The second flow configuration investigated is the incompressible shear-driven cavity flow. It is a geometrically-induced separated boundary layer flow having a number of applications in aeronautics. The leading two-dimensional instability of the flow is mostly localised along the shear layer developing at the interface between the outer boundary layer flow and the inner cavity flow (Sipp et al., 2010). This oscillatory global instability of the external shear layer relies on two essential mechanisms. On the one hand, the convectively unstable nature of the shear layer causes

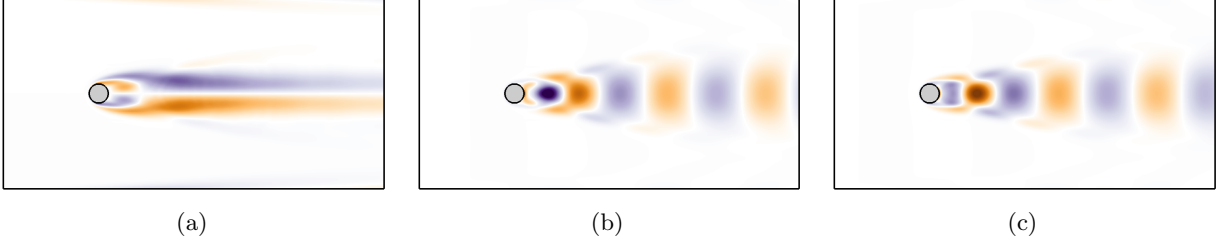


Figure 3: Vorticity fields of (a) the shift mode, (b) the first and (c) second POD modes of the cylinder flow at  $Re = 100$ .

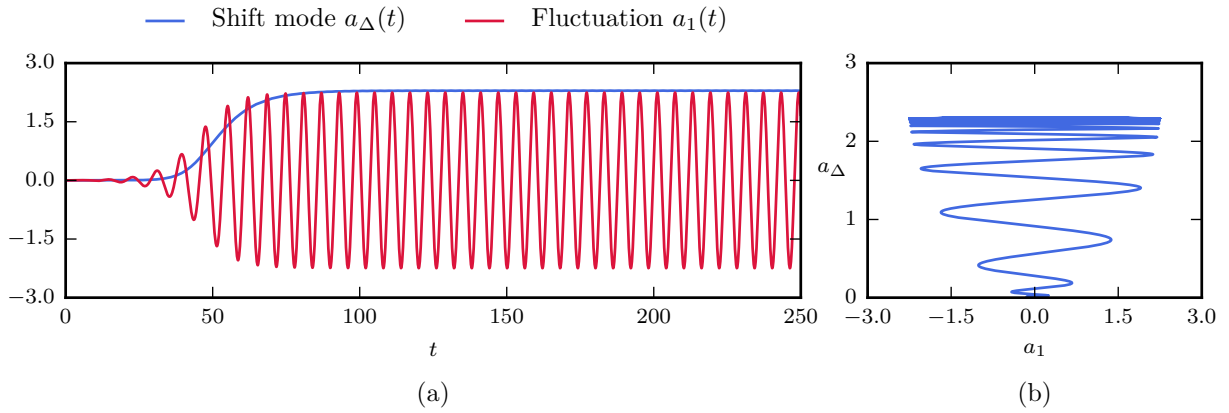


Figure 4: (a) Time evolution of the POD coefficients for the cylinder flow at  $Re = 100$ . The time evolution of  $a_2(t)$ , not shown, is very similar to that of  $a_1(t)$ . (b) Trajectory in the phase space projected onto the  $a_1 - a_\Delta$  plane.

perturbations to grow as they are convected downstream, while on the other hand, the feedback mechanism provided by the inner-cavity recirculating flow allows these same perturbations to eventually re-excite the upstream shear layer. The coupling between these two mechanisms gives rise to a linearly unstable feedback loop at sufficiently high Reynolds numbers. Note that for compressible shear-driven cavity flows, a similar unstable feedback loop exists wherein the feedback mechanism is provided by upstream-propagating acoustic waves (Rossiter, 1964; Rowley et al., 2002; Yamouni et al., 2013). This strictly two-dimensional linearly unstable flow configuration has served multiple purposes over the past decade: illustration of optimal control and reduced-order modelling (Barbagallo et al., 2009), investigation of the nonlinear saturation process of globally unstable flows (Sipp and Lebedev, 2007), or as an introduction to dynamic mode decomposition (Schmid, 2010), to name just a few.

The computational domain and boundary conditions considered are the same as in Sipp and Lebedev (2007). The Reynolds number is set to  $Re = 4250$ , based on the free-stream velocity  $U_\infty$  and the depth  $L$  of the open cavity. As for the cylinder, the linearly unstable flow, the corresponding eigenspectrum and the vorticity field of the leading unstable eigenmode are presented in figure 5 for the sake of completeness. Using this linearly unstable flow as initial condition, another direct numerical simulation has been run until a statistically steady state is achieved. The vorticity field of the corresponding shift mode and of the first dominant POD mode are shown in figure 6(a) and (b), respectively. While the leading unstable eigenmode and the dominant POD mode of the cylinder flow are extremely different, this is not the case for the shear-driven cavity flow at  $Re = 4250$ . Comparing figure 5(c) and figure 6(b), it can be seen that these two modes are now very similar. The evolution in time of the coefficients  $a_1(t)$  (dominant POD mode) and  $a_\Delta(t)$  (shift mode) is shown in figure 7. Note that these curves appear as filled-in regions due to the high-frequency oscillations of  $a_1(t)$ . Despite the fundamental difference of the geometry, the different frequency of the oscillations and the smaller growth rate of the instability, the two flows considered herein appear to exhibit relatively similar dynamics when looking at the systems' trajectories projected onto the  $a_1$ - $a_\Delta$  planes: both low-dimensional representations of the flows appear to evolve along a parabolic manifold, see figure 4(b) and figure 7(b).

## 5 Results and discussion

Following the seminal work of Noack et al. (2003), so-called quadratic *Galerkin Regression models* are made of the basic building blocks necessary for reduced-order modelling of the flow configurations considered, *i.e.* a linear operator  $\tilde{\mathcal{L}}$  and an energy-preserving quadratic nonlinearity  $\tilde{\mathcal{N}}(\mathbf{a})$ . For that purpose, the library  $\Theta(\mathbf{a})$  used in the identification process is defined as  $P_2(\mathbf{a})$ , *i.e.* all the polynomials of degree 2 or less in the entries of  $\mathbf{a}$ . The quadratic Galerkin regression models identified for the cylinder flow and the shear-driven cavity flow are reported in tables 1 and 2. A second type of models, cubic Galerkin Regression models, are made of the same basic building blocks as their quadratic counterparts. They moreover include higher-order nonlinearities which can serve to model the truncated modes, as discussed in §3.2. For that purpose, the library  $\Theta(\mathbf{a})$  used in the identification process is defined as  $P_3(\mathbf{a})$ , *i.e.* all the polynomials of degree 3 or less in the entries of  $\mathbf{a}$ . Up to 57 coefficients then need to be identified for the present case with  $n = 3$  state variables. The cubic models identified for the cylinder flow and the shear-driven cavity flow are reported in tables 3 and 4.

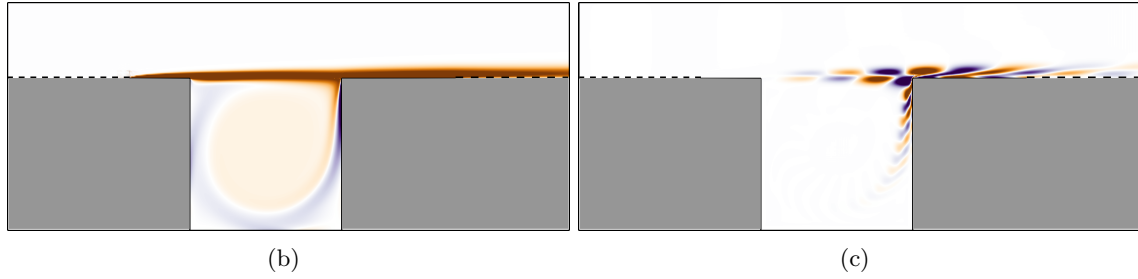
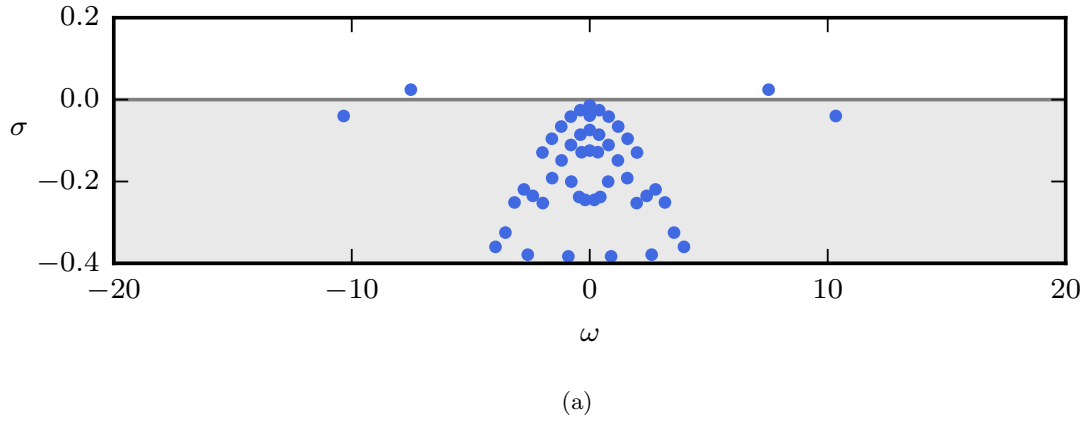


Figure 5: (a) Eigenspectrum of the linearised Navier-Stokes operator for the shear-driven cavity flow at  $Re = 4250$ . Vorticity fields of (b) the base flow and (c) the leading linearly unstable eigenmode. The dashed lines indicate the spatial extent over which the free-slip boundary condition is imposed.

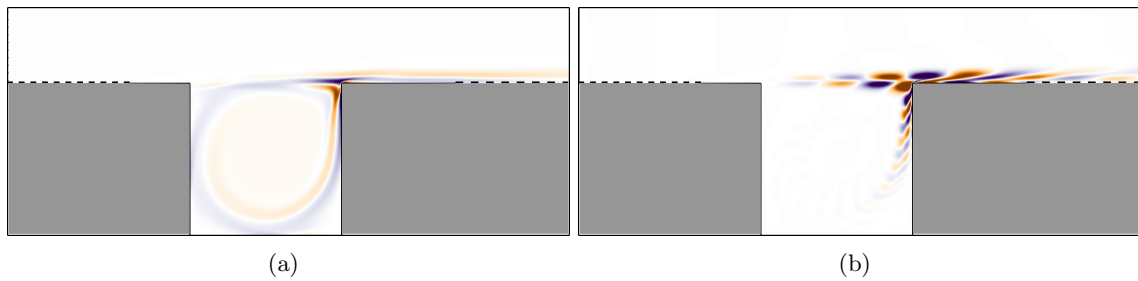


Figure 6: Vorticity fields of (a) the shift mode and (b) the first POD mode for the shear-driven cavity flow at  $Re = 4250$ .

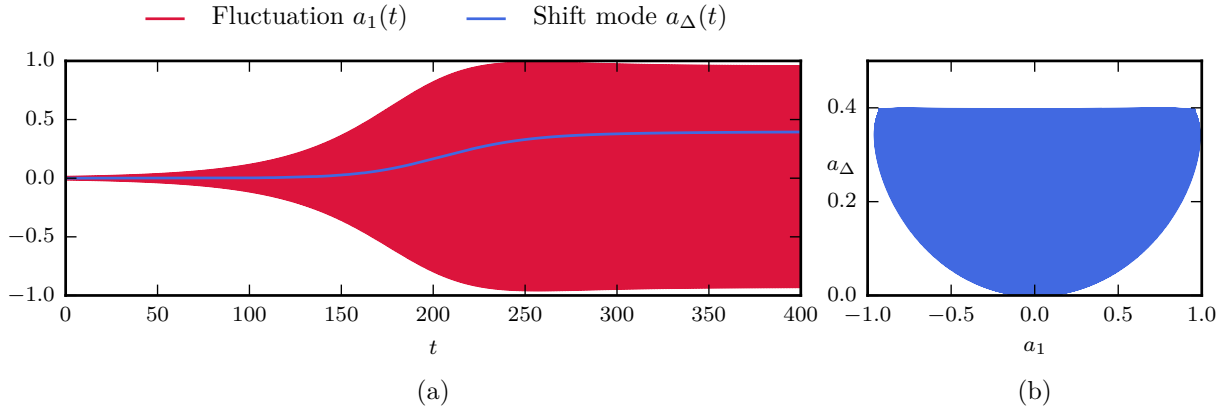


Figure 7: (a) Time evolution of the POD coefficients for the shear-driven cavity flow at  $Re = 4250$ . The time evolution of  $a_2(t)$ , not shown, is very similar to that of  $a_1(t)$ . (b) Trajectory in the phase space projected onto the  $a_1 - a_\Delta$  plane. Note that, for both figures, the  $a_i(t)$  coefficients have been multiplied by 100.

## 5.1 Cylinder flow

Figures 8 and 9 provide a comparison of the dynamics predicted by the low-dimensional Galerkin Regression models identified using constrained sparse regression against the dynamics of the original system for the two-dimensional cylinder flow at  $Re = 100$ . It also provides the dynamics predicted by two additional data-driven reduced-order models, namely:

- the minimal Galerkin projection model including only the shift mode and the first two POD modes,
- a Galerkin projection model including the shift mode and the first eight POD modes.

Figure 8 depicts the evolution of the mean flow distortion as a function of time for the different reduced-order models. As reported in previous works (Noack et al., 2003; Rowley and Dawson, 2016), the low-dimensional systems derived based on a Galerkin projection procedure that includes only the shift mode and the leading POD modes significantly over-estimate the duration of the transients. As explained by Noack et al. (2003), this over-estimation results from the fact that the leading POD modes (see figure 3) provide only a crude approximation of the leading linear instability eigenmodes (see figure 2). These Galerkin projection models moreover suffer from an energy overshoot once nonlinear saturation kicks in. This overshoot and the ensuing larger amplitude of the mean flow distortion mostly result from the disruption of the energy cascade due to neglecting the higher-harmonic POD modes. Being neglected, these higher harmonics cannot absorb the excess energy produced by the two most energetic modes. The latter then grow beyond the correct value until the mean-flow distortion  $a_\Delta(t)$  can eventually absorb this excess energy via the coupling terms. As shown in figure 8(b), the quadratic *Galerkin Regression* model suffers from similar drawbacks, although the duration of transients is shortened and the final amplitude of the mean flow distortion is in agreement with that of the original system.

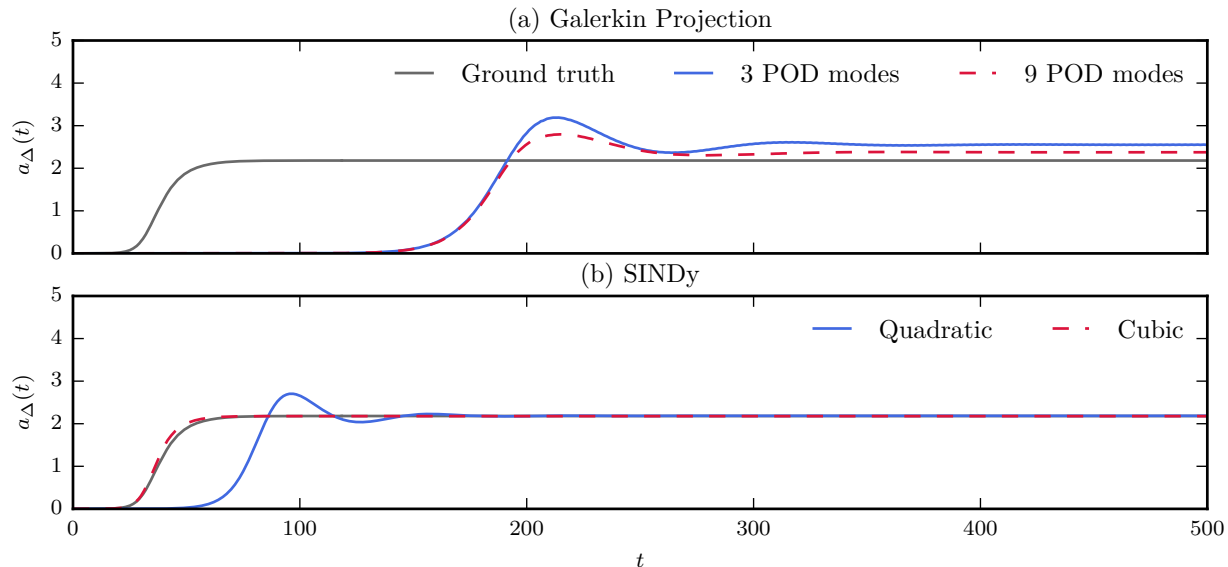


Figure 8: Comparison of the time-evolution of the mean flow distortion  $a_{\Delta}$  predicted by the different data-driven models for the two-dimensional cylinder flow at  $Re = 100$ .

The dynamics predicted by the cubic Galerkin Regression model are shown in figures 8(b) figure 9(d). It can be seen that including higher-order nonlinearities results in a cubic Galerkin regression model that provides an almost perfect fit to the original data. The amplitude of the limit cycle is less than 0.5% higher than that of the original system while the saturation of the mean flow distortion differs by less than 0.1%. It has to be noted however that the growth rate of the instability is slightly over-estimated. Nonetheless, the inclusion of the cubic nonlinearities has a stabilising effect, hence preventing the energy overshoot and/or larger limit cycle amplitude observed for the quadratic models. A stabilising cubic term would also be obtained when deriving the Landau amplitude equation describing the transient dynamics of a small perturbation in the neighborhood of a supercritical Hopf bifurcation (Noack and Eckelmann, 1994; Sipp and Lebedev, 2007). The excellent predictions of the cubic model as well as the existing connections with amplitude equations (Noack and Eckelmann, 1994; Sipp and Lebedev, 2007), adiabatic elimination (Haken, 1983) or center manifold reduction (Carini et al., 2015; Wiggins, 2003) thus justify *a posteriori* the use of higher-order nonlinearities to model the influence of the truncated modes onto the driving ones as discussed in §3.2.

## 5.2 Shear-driven cavity flow

As for the cylinder flow, figures 10 and 11 provide a comparison of the dynamics predicted by the low-dimensional Galerkin Regression models identified using constrained sparse regression against the dynamics of the original system for the shear-driven cavity flow at  $Re = 4250$ . The dynamics predicted by unconstrained SINDy models have also been included along with those predicted by the corresponding Galerkin projection models. Note that to identify meaningful Galerkin regression models, the sparse regression algorithm requires a pre-processing step so that all the features in  $\Theta(\mathbf{a})$  have the same range in order to facilitate the optimisation procedure. Although the geometry



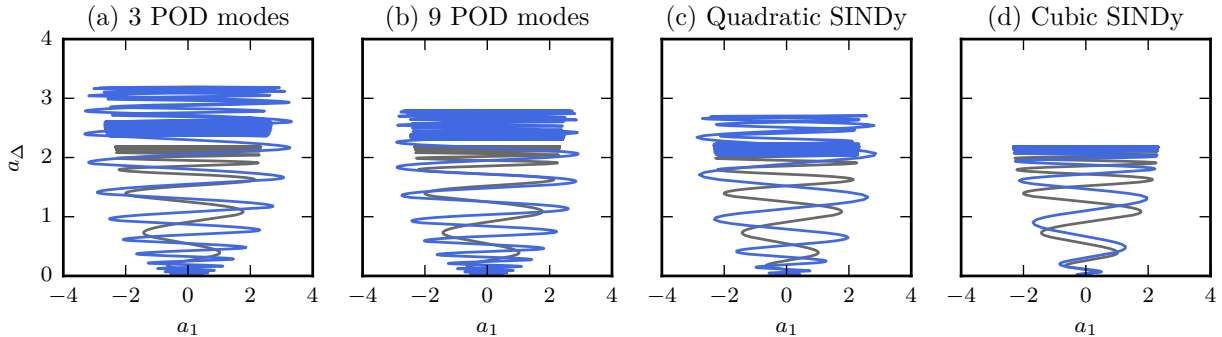


Figure 9: Comparison of the trajectory in the  $a_1 - a_\Delta$  plane predicted by the different reduced-order models for the two-dimensional cylinder flow at  $Re = 100$ . The light gray trajectory is the one given by a direct numerical simulation.

and the physics are quite different from that of the two-dimensional cylinder flow, it can be seen that the present *Galerkin projection* models suffer from similar drawbacks as before: a misprediction of the transients duration and the saturation to higher mean flow distortion due to the disruption of the energy cascade. However, the key difference is that for the shear-driven cavity flow, the growth rate of the linear instability mode is slightly over-predicted by the Galerkin projection models.

Let us now draw our attention onto the quadratic Galerkin Regression models. Looking at the second subplot of figures 10 and 11, it can be seen that both models correctly reproduce the asymptotic dynamics of the shear-driven cavity flow. The major difference however relies in the prediction of the transient dynamics. Although it would appear as the most physical one, the quadratic Galerkin regression model with an energy-preserving quadratic nonlinearity severely over-predicts the duration of the transients. Comparatively, the unconstrained quadratic model yields a much better prediction despite the small unphysical overshoot observed at the onset of nonlinear saturation. It is not clear at the present time why the constrained quadratic model performs so badly. One way to improve its performance would be to constrain the eigenspectrum of the low-dimensional linear operator to be a subset of its high-dimensional counterpart. Such a constraint on the determinant of the low-dimensional linear operator is however a non-convex constraint and does not fall in the scope of the library CVXOPT used in the present work.

Finally, it can be seen in figures 10(c) and 11(c) that both cubic models exhibit similar accuracy. The only visible difference between these two models is that the growth rate of the linear instability is slightly over-estimated by the constrained model while being slightly under-estimated by the unconstrained one. Given the similar performance, one might thus wonder what is the benefit of constraining the identification process. The answer to this question relies in the eigenspectrum of the low-dimensional linear operator. For the unconstrained model, this matrix and its eigenspectrum are given by

$$\tilde{\mathcal{L}} = \begin{bmatrix} 0 & 8.011 & 0.0408 \\ -7.0579 & 0.0465 & -0.1146 \\ -0.0181 & 0 & 0.0191 \end{bmatrix} \text{ and } \Lambda = \begin{bmatrix} 0.0231 + i7.519 & 0 & 0 \\ 0 & 0.0231 - i7.519 & 0 \\ 0 & 0 & 0.0194 \end{bmatrix}.$$

The unconstrained model hence correctly identifies the fixed point of the system as being a linearly unstable spiral within the  $a_1 - a_2$  plane. It also identifies it as being linearly unstable in the  $a_\Delta$  direction. Naively, this result appears consistent. Looking at the time-evolution depicted in figure 10(c) without prior knowledge of the problem, one could easily conclude that the system is linearly unstable in the  $a_\Delta$  direction. From an identification point of view, the governing equations for  $a_1$ ,  $a_2$  and  $a_\Delta$  are obtained independently from one another in the absence of constraints that would otherwise couple them. As a consequence, an equation predicting a linear instability of  $a_\Delta$  is thus the simplest model identifiable which balances parsimony and consistency with observed measurements. However, given our prior knowledge about the physics of the problem, this is not an acceptable model. It could lead to a misunderstanding of the physics at play and seriously alter the practical performance of a linear or nonlinear controller based on such a faulty reduced-order model. As a comparison, the low-dimensional linear operator of the constrained cubic Galerkin regression model and the corresponding eigenspectrum are given by

$$\tilde{\mathcal{L}} = \begin{bmatrix} 0.0063 & 8.2337 & 0.1442 \\ -7.2843 & 0.049 & 0 \\ -0.0347 & 0 & -0.0243 \end{bmatrix} \text{ and } \Lambda = \begin{bmatrix} 0.0276 + i7.75 & 0 & 0 \\ 0 & 0.0276 - i7.75 & 0 \\ 0 & 0 & -0.0243 \end{bmatrix}.$$

Given the linearly stable nature of the  $a_\Delta$  direction now predicted, it is clear that coupling all of the equations governing the evolution of the system through the use of constraints mimicking the energy-preserving nature of the quadratic nonlinearity enables the identification of a much more physical low-dimensional system.

## 6 Conclusion

This paper develops a new data-driven *Galerkin regression* framework to identify nonlinear reduced-order models of a fluid. The resulting models incorporate a number of beneficial features of standard Galerkin projection, making them easy to interpret and use, but without the need for access to a high-fidelity Navier-Stokes model for the projection. Galerkin regression models also provide a more flexible model identification, in that they readily generalize to include higher-order nonlinear terms that model the effect of truncated modes; the inclusion of these terms is shown to be extremely effective in the examples presented here. The Galerkin regression framework leverages the recent sparse identification of nonlinear dynamics (SINDy) algorithm (Brunton et al., 2016b), and significantly generalizes it to include user-provided constraints directly into the sparsity-promoting regression. These additional constraints can be used to enforce *a priori* known values of some of the regression coefficients, inherent symmetries of the system of equations or some physical behaviour such as the energy-preserving nature of the quadratic nonlinearity of the Navier-Stokes equations.

The Lorenz system, the two-dimensional cylinder flow and the shear-driven cavity have each been carefully analyzed to illustrate the system identification capabilities of the resulting algorithm. For that purpose, two polynomial libraries have been used and the constraints have been chosen in order to enforce different physical properties. The accuracy and performance of the so-called *Galerkin regression* models have been compared against reduced-order models derived using a classical Galerkin projection method. All of the regression models qualitatively reproduce the main features of the original system: linear instability of the fixed point and final saturation to a periodic

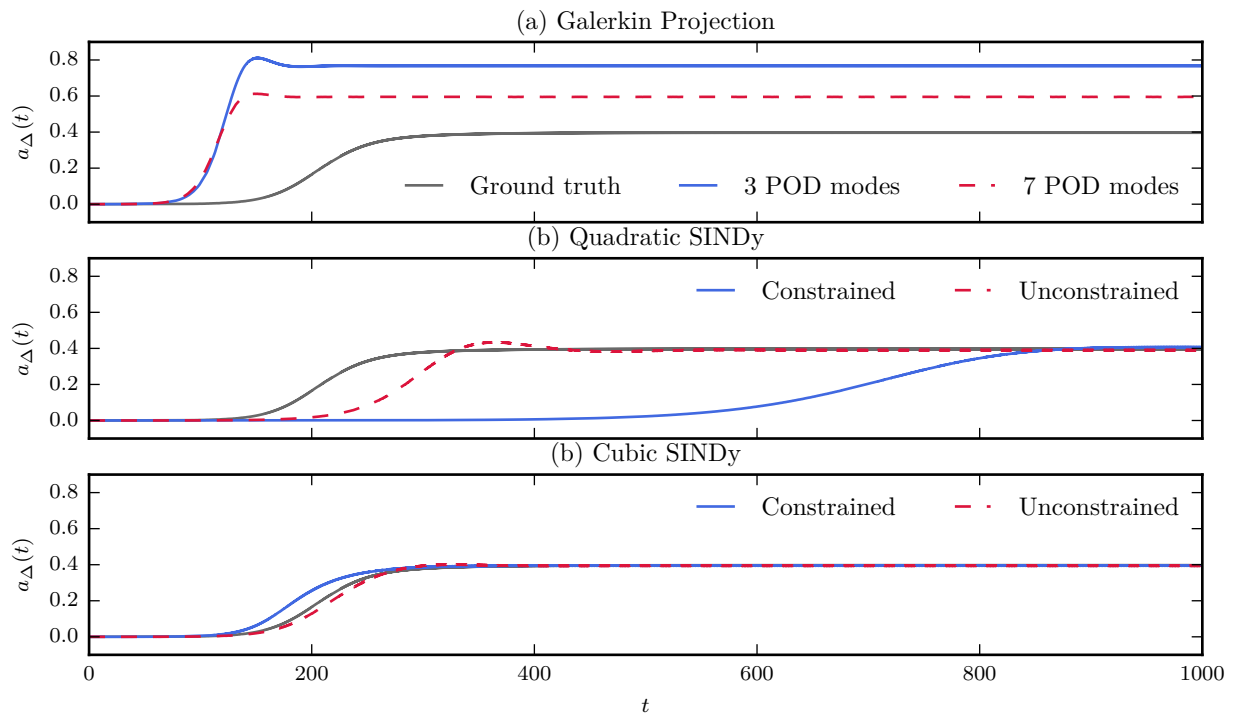


Figure 10: Comparison of the time-evolution of the mean flow distortion  $a_{\Delta}$  predicted by the different data-driven models for the two-dimensional shear-driven cavity flow at  $Re = 4250$ .

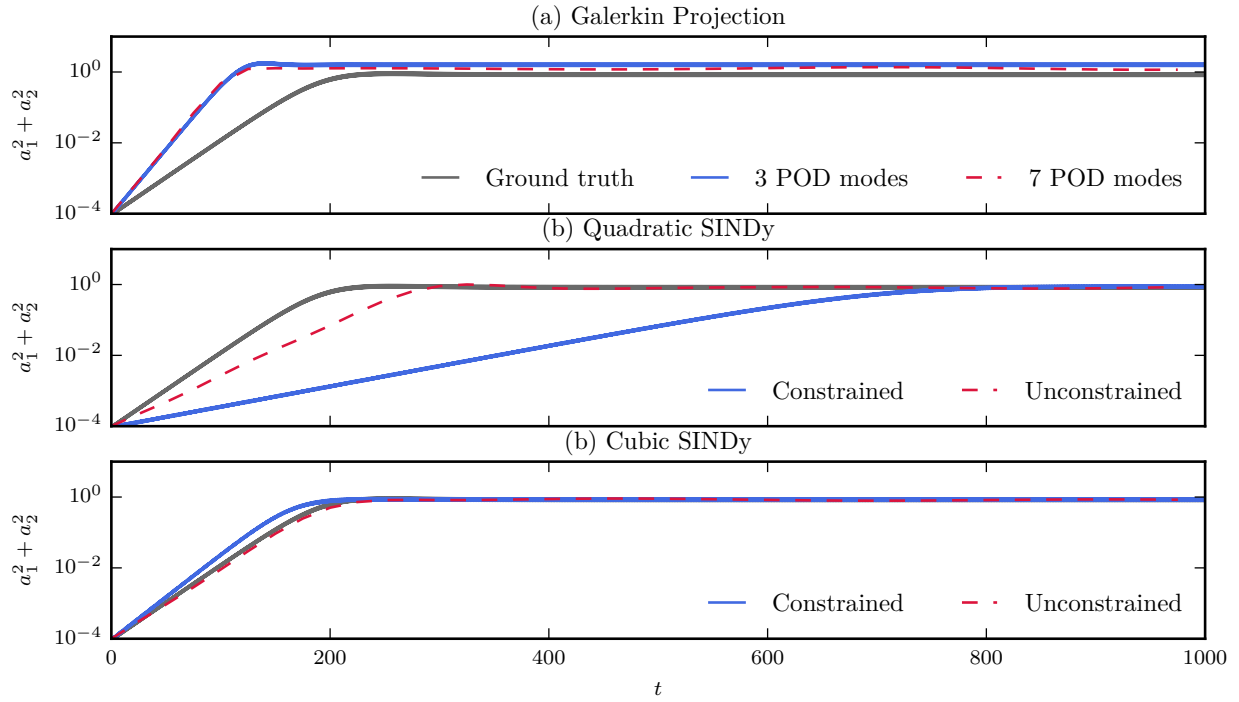


Figure 11: Comparison of the time-evolution of the fluctuation's kinetic energy  $a_1^2 + a_2^2$  predicted by the different data-driven models for the two-dimensional shear-driven cavity flow at  $Re = 4250$ .

limit cycle. Though these models rely essentially on a data-driven approach, visual inspection of their trajectories in the phase space highlights the connection between the quadratic models and the models obtained using a Galerkin projection procedure in the seminal work of Noack et al. (2003). Moreover, both flow configurations highlight the importance of including cubic nonlinearities into the admissible pool of functions for the identification process, something utterly impossible with classical Galerkin projection without significant additional post-analysis. These cubic terms then model the influence of the truncated modes onto the driving ones, eventually enabling the identification of a low-dimensional system with much better predictive capabilities. Although the unconstrained cubic low-dimensional model of the shear-driven cavity reproduces faithfully the dynamics of the original system, this particular flow configuration has highlighted the importance of incorporating physically meaningful constraints into the regression to ensure that the identified model has the correct physical behaviour. In their absence, the SINDy algorithm incorrectly identifies the mean flow distortion as a linearly unstable manifold of the fixed point, while adding constraints results in the correct identification of a linearly stable eigenvalue.

Despite its promise, such an approach to system identification still suffers from certain limitations. One such limitation is illustrated by the quadratic constrained model identified for the shear-driven cavity flow which strongly under-estimates the growth rate of the linear instability. Given prior knowledge of the linear stability of the high-dimensional system, see §4.2, one could then constrain the eigenspectrum of the low-dimensional linear operator to be a subset of its high-dimensional counterpart. Such a constraint, involving the determinant of the low-dimensional matrix, falls outside the scope of convex optimisation. Current developments, based on the nonlinear optimisation library NLOPT (Johnson, 2014), attempt to overcome such limitations. One might also argue that the systems considered in the present work are inherently low-dimensional and are thus not representative of the high-dimensionality of a transitional or turbulent flow. However, such flows have already been modelled with some success using a Galerkin projection procedure (Gloerfelt, 2008). Given the parallels drawn in the present work between Galerkin projection and Galerkin regression, there is reason to believe that the present approach may be successfully applied to such flows as well. Indeed, this is an exciting future direction and is the subject of ongoing work. Including high-order nonlinear terms in the pool of admissible functions in combination with the sparsity-promoting capabilities of the algorithm might furthermore allow the identification of smaller and more robust reduced-order models without significantly altering their accuracy and predictive capabilities.

## Acknowledgment

We are grateful for many fruitful discussions with Bernd Noack, Josh Proctor and Nathan Kutz. We also appreciate valuable feedback from Scott Dawson and Clancy Rowley. SLB acknowledges generous funding support from the Defense Advanced Research Projects Agency (DARPA HR0011-16-C-0016) and from the Air Force Office of Scientific Research (AFOSR FA9550-13-1-0183).

## A Coefficients of the different models identified

The following tables provide the coefficients for each model identified using the SINDy algorithm extended with the energy-preserving constraint for the quadratic nonlinear term. Models A1 (see

	$\dot{a}_1$	$\dot{a}_2$	$\dot{a}_\Delta$
$a_1$	0.0523	0.6667	0
$a_2$	-0.6856	0.0617	0
$a_\Delta$	0	0	-0.0513
$a_1^2$	0	0	0.0245
$a_1 a_2$	0	0	0
$a_1 a_\Delta$	-0.0245	0.1599	0
$a_2^2$	0	0	0.025
$a_2 a_\Delta$	-0.1599	-0.025	0
$a_\Delta^2$	0	0	0

Table 1: Coefficients of the quadratic Galerkin regression model for the two-dimensional cylinder flow at  $Re = 100$ .

	$\dot{a}_1$	$\dot{a}_2$	$\dot{a}_\Delta$
$a_1$	$5.092 \cdot 10^{-3}$	-7.068	0
$a_2$	7.987	$8.166 \cdot 10^{-3}$	0
$a_\Delta$	$-2.369 \cdot 10^{-2}$	0.219	-0.034
$a_1^2$	0	-0.1543	0
$a_1 a_2$	0.1542	0.4106	0
$a_1 a_\Delta$	0	-0.0527	0
$a_2^2$	-0.4106	0	0.0343
$a_2 a_\Delta$	0.0527	-0.0343	0
$a_\Delta^2$	0	0	0

Table 2: Coefficients of the quadratic Galerkin regression model A2 for the two-dimensional shear-driven cavity flow at  $Re = 4250$ .

table 1) and B1 (see table 3) are the quadratic and cubic Galerkin regression models obtained for the two-dimensional cylinder flow at  $Re = 100$ , respectively. Their counterparts for the shear-driven cavity flow, *i.e.* models A2 and B2, are given in tables 2 and 4.

## B Influence of the sparsity knob $\lambda$ and model selection

Although it has not been discussed in the core of the present paper, the choice of the sparsity knob  $\lambda$  is of crucial importance in the selection of the final model. Governing the level of sparsity, this parameter  $\lambda$  is thus directly related to the accuracy and complexity of the identified models. If  $\lambda$  is too small, very few terms are eliminated and the identified model has an artificially high complexity. On the other hand, if  $\lambda$  is too large, the identified model may have too few terms, thus impacting its accuracy. To evaluate *a priori* the predictive capabilities of the identified model, it is convenient to analyze the number of non-zero coefficients and the  $r^2$  score (Draper and Smith, 2014) (also known as the coefficient of determination) as a function of the sparsity knob  $\lambda$ .

Figure 12 depicts the evolution of these two metrics as a function of the sparsity knob  $\lambda$  for the cubic Galerkin regression model of the cylinder flow. It can be observe that increasing  $\lambda$  up to almost 1 has a negligible influence of the *a priori* accuracy of the equations identified for the

	$\dot{a}_1$	$\dot{a}_2$	$\dot{a}_\Delta$
$a_1$	0.0768	0.7527	0
$a_2$	-0.745	0.1046	0
$a_\Delta$	0	0	-0.0357
$a_1^2$	0	0	0.0596
$a_1 a_2$	0	0	0
$a_1 a_\Delta$	-0.0596	0.1237	0
$a_2^2$	0	0	0.0641
$a_2 a_\Delta$	-0.1236	-0.0641	0
$a_\Delta^2$	0	0	0
$a_1^3$	0	-0.0264	0
$a_1^2 a_2$	0.0318	-0.005	0
$a_1^2 a_\Delta$	0	0	-0.0189
$a_1 a_2^2$	0	-0.0275	0
$a_1 a_2 a_\Delta$	0	0	0
$a_1 a_\Delta^2$	0.0107	0.025	0
$a_2^3$	0.0323	-0.005	0
$a_2^2 a_\Delta$	0	0	-0.0208
$a_2 a_\Delta^2$	-0.0358	0.0135	0
$a_\Delta^3$	0	0	0

Table 3: Coefficients of the cubic Galerkin regression model B1 for the two-dimensional cylinder flow at  $Re = 100$ .

	$\dot{a}_1$	$\dot{a}_2$	$\dot{a}_\Delta$
$a_1$	$6.33 \cdot 10^{-3}$	-7.284	-0.0347
$a_2$	8.233	0.049	0
$a_\Delta$	0.144	0	-0.0243
$a_1^2$	0	-0.0765	0
$a_1 a_2$	0.0765	0.088	0
$a_1 a_\Delta$	0	13.53	1.139
$a_2^2$	-0.0881	0	0.0364
$a_2 a_\Delta$	-13.53	-0.036	-0.351
$a_\Delta^2$	-1.1393	0.351	0
$a_1^3$	0.0053	-2.219	-0.1659
$a_1^2 a_2$	2.176	-0.063	0.063
$a_1^2 a_\Delta$	0.0805	0.549	0
$a_1 a_2^2$	0	-2.487	-0.184
$a_1 a_2 a_\Delta$	0.2242	-0.389	0
$a_1 a_\Delta^2$	-0.0398	-20.499	-1.725
$a_2^3$	2.445	-0.0592	0.0723
$a_2^2 a_\Delta$	0	0.0602	-0.0195
$a_2 a_\Delta^2$	20.472	0.0579	0.527
$a_\Delta^3$	1.805	-2.75	0

Table 4: Coefficients of the cubic Galerkin regression model B2 for the two-dimensional shear-driven cavity flow at  $Re = 4250$ .

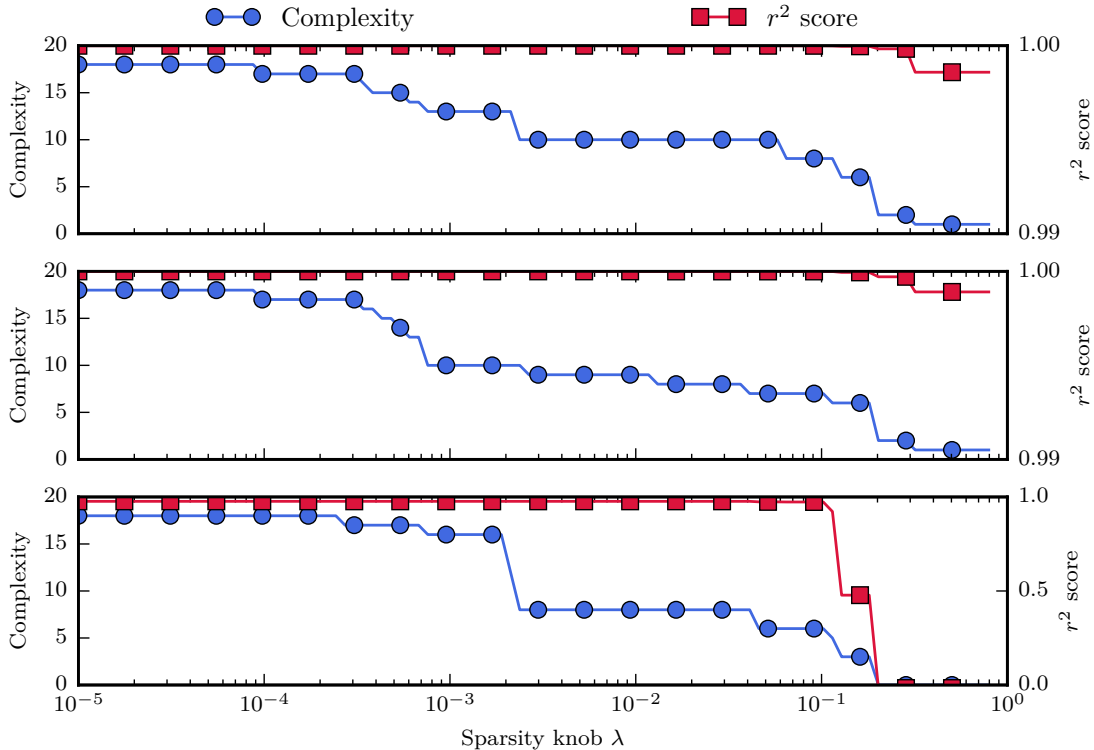


Figure 12: Illustration of the influence of the sparsity knob  $\lambda$  on the number of coefficients retained in the cubic SINDy model for the two-dimensional cylinder flow at  $Re = 100$ . Evolution of the number of non-zero coefficients and of the  $r^2$  score as a function of  $\lambda$  for the governing equation of  $a_1$  (top),  $a_2$  (middle) and  $a_\Delta$  (bottom).

fluctuation’s dynamics. Paradoxically, all the coefficients for the mean flow distortion’s governing equation are set to zero for  $\lambda > 0.1$ , thus highlighting the over-aggressive sparsity promotion of the algorithm for this particular knob. This sudden drop in the  $r^2$  score and number of non-zero coefficients of the mean flow distortion equation corresponds to the existence of a kink in the Pareto fronts of all three equations. The corresponding model is the one that provides the highest *a priori* accuracy while having the lowest complexity. All the identified models presented in this work have been selected using the same criterion. It has to be noted that such a model selection strategy can be combined with more sophisticated K-fold cross-validation to get an even better estimate of the *a priori* accuracy of the identified models.

## References

- E. Åkervik, L. Brandt, D. S. Henningson, J. Höpfner, O. Marxen, and P. Schlatter. Steady solutions of the navier-stokes equations by selective frequency damping. *Physics of Fluids (1994-present)*, 18(6):068102, 2006.
- M. S. Andersen, J. Dahl, and L. Vandenberghe. CVXOPT: A Python package for convex optimiza-



- tion, Version 1.1.6, 2013.
- S. Bagheri. Koopman-mode decomposition of the cylinder wake. *J. Fluid Mech*, 726:596–623, 2013.
- S. Bagheri, L. Brandt, and D.S. Henningson. Input-output analysis, model reduction and control of the flat-plate boundary layer. *J. Fluid Mech.*, 620:263–298, 2009.
- M. J. Balajewicz, E. H. Dowell, and B. R. Noack. Low-dimensional modelling of high-Reynolds-number shear flows incorporating constraints from the Navier–Stokes equation. *J. Fluid Mech.*, 729:285–308, 2013.
- A. Barbagallo, D. Sipp, and P. J. Schmid. Closed-loop control of an open cavity flow using reduced-order models. *Journal of Fluid Mechanics*, 641:1–50, 2009.
- D. Barkley. Linear analysis of the cylinder wake mean flow. *EPL (Europhysics Letters)*, 75(5):750, 2006.
- D. Barkley and R. D. Henderson. Three-dimensional Floquet stability analysis of the wake of a circular cylinder. *J. Fluid Mech.*, 322:215–241, 1996.
- G. Berkooz, P. J. Holmes, and J. L. Lumley. The proper orthogonal decomposition in the analysis of turbulent flows. *Annual review of fluid mechanics*, 25(1):539–575, 1993.
- J. Bongard and H. Lipson. Automated reverse engineering of nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*, 104(24):9943–9948, 2007.
- S. L. Brunton and B. R. Noack. Closed-loop turbulence control: Progress and challenges. *Applied Mechanics Reviews*, 67(5):050801, 2015.
- S. L. Brunton, J. L. Proctor, and J. N. Kutz. Sparse identification of nonlinear dynamics with control (SINDYc). *arXiv preprint arXiv:1605.06682*, 2016a.
- S. L. Brunton, J. L. Proctor, and J. N. Kutz. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*, 113(15):3932–3937, 2016b.
- E. J. Candès. Compressive sensing. *Proceedings of the International Congress of Mathematics*, 2006.
- M. Carini, F. Auteri, and F. Giannetti. Centre-manifold reduction of bifurcating flows. *J. Fluid Mech.*, 767:109–145, 3 2015. ISSN 1469-7645. doi: 10.1017/jfm.2015.3. URL [http://journals.cambridge.org/article\\_S0022112015000038](http://journals.cambridge.org/article_S0022112015000038).
- K. Carlberg, R. Tuminaro, and P. Boggs. Preserving lagrangian structure in nonlinear model reduction with application to structural dynamics. *SIAM Journal on Scientific Computing*, 37(2):B153–B184, 2015.
- R. Chartrand. Numerical differentiation of noisy, nonsmooth data. *ISRN Applied Mathematics*, 2011, 2011.
- D. L. Donoho. Compressed sensing. *IEEE Transactions on Information Theory*, 52(4):1289–1306, 2006.

- N. R. Draper and H. Smith. *Applied regression analysis*. John Wiley & Sons, 2014.
- N. Fabbiane, O. Semeraro, S. Bagheri, and D. S. Henningson. Adaptive and model-based control theory applied to convectively unstable flows. *Applied Mechanics Reviews*, 66(6):060801, 2014.
- P.F. Fischer, J.W. Lottes, and S.G. Kerkemeir. Nek5000 Web pages, 2008. <http://nek5000.mcs.anl.gov>.
- X. Gloerfelt. Compressible proper orthogonal decomposition/galerkin reduced-order model of self-sustained oscillations in a cavity. *Physics of Fluids*, 20(11):115105, 2008.
- G. H. Golub and C. F. Van Loan. *Matrix computations*, volume 3. JHU Press, 2012.
- H. Haken. Springer series in synergetics. *Editors: M. Cardona P. Fulde H.-J. Queisser*, page 269, 1983.
- P. J. Holmes, J. L. Lumley, G. Berkooz, and C. W. Rowley. *Turbulence, coherent structures, dynamical systems and symmetry*. Cambridge Monographs in Mechanics. Cambridge University Press, Cambridge, England, 2nd edition, 2012.
- M. Ilak and C. W. Rowley. Modeling of transitional channel flow using balanced proper orthogonal decomposition. *Physics of Fluids*, 20:034103, 2008.
- S. J. Illingworth, A. S. Morgans, and C. W. Rowley. Feedback control of flow resonances using balanced reduced-order models. *Journal of Sound and Vibration*, 330(8):1567–1581, 2010.
- S. G. Johnson. The nlopt nonlinear-optimization package, 2014.
- J.-N. Juang and R. S. Pappa. An eigensystem realization algorithm for modal parameter identification and model reduction. *Journal of guidance, control, and dynamics*, 8(5):620–627, 1985.
- E. Kaiser, B. R. Noack, L. Cordier, A. Spohn, M. Segond, M. Abel, G. Daviller, J. Osth, S. Krajinovic, and R. K. Niven. Cluster-based reduced-order modelling of a mixing layer. *J. Fluid Mech.*, 754:365–414, 2014.
- J. N. Kutz, S. L. Brunton, B. W. Brunton, and J. L. Proctor. *Dynamic Mode Decomposition: Data-Driven Modeling of Complex Systems*. SIAM, 2016.
- E. N. Lorenz. Deterministic nonperiodic flow. *Journal of the atmospheric sciences*, 20(2):130–141, 1963.
- A. J. Majda and J. Harlim. Physics constrained nonlinear regression models for time series. *Nonlinearity*, 26(1):201, 2012.
- N. M. Mangan, S. L. Brunton, J. L. Proctor, and J. N. Kutz. Inferring biological networks by sparse identification of nonlinear dynamics. *To appear in the IEEE Transactions on Molecular, Biological, and Multi-Scale Communications*, 2016.
- I. Mezić. Spectral properties of dynamical systems, model reduction and decompositions. *Nonlinear Dynamics*, 41(1-3):309–325, 2005.

- I. Mezić. Analysis of fluid flows via spectral properties of the Koopman operator. *Annual Review of Fluid Mechanics*, 45:357–378, 2013.
- A. G. Nair and K. Taira. Network-theoretic approach to sparsified discrete vortex dynamics. *Journal of Fluid Mechanics*, 768:549–571, 2015.
- B. R. Noack and H. Eckelmann. A global stability analysis of the steady and periodic cylinder wake. *J. Fluid Mech.*, 270:297–330, 1994.
- B. R. Noack, K. Afanasiev, M. Morzynski, G. Tadmor, and F. Thiele. A hierarchy of low-dimensional models for the transient and post-transient cylinder wake. *J. Fluid Mech.*, 497:335–363, 2003.
- B. R. Noack, W. Stankiewicz, M. Morzynski, and P. J. Schmid. Recursive dynamic mode decomposition of a transient cylinder wake. *arXiv preprint arXiv:1511.06876*, 2015.
- J. E. Rossiter. Wind tunnel experiments on the flow over rectangular cavities at subsonic and transonic speeds. Technical report, Ministry of Aviation; Royal Aircraft Establishment; RAE Farnborough, 1964.
- C. W. Rowley. Model reduction for fluids using balanced proper orthogonal decomposition. *International Journal of Bifurcation and Chaos*, 15(3):997–1013, 2005.
- C. W. Rowley and S. Dawson. Model reduction for flow analysis and control. *Annual Review of Fluid Mechanics*, 49(1), 2016.
- C. W. Rowley, T. Colonius, and A. J. Basu. On self-sustained oscillations in two-dimensional compressible flow over rectangular cavities. *Journal of Fluid Mechanics*, 455:315–346, 2002.
- C. W. Rowley, I. Mezić, S. Bagheri, P. Schlatter, and D.S. Henningson. Spectral analysis of nonlinear flows. *J. Fluid Mech.*, 645:115–127, 2009.
- S. H. Rudy, S. L. Brunton, J. L. Proctor, and J. N. Kutz. Data-driven discovery of partial differential equations. *arXiv preprint arXiv:1609.06401*, 2016.
- P. J. Schmid. Dynamic mode decomposition of numerical and experimental data. *Journal of Fluid Mechanics*, 656:5–28, 2010.
- M. Schmidt and H. Lipson. Distilling free-form natural laws from experimental data. *science*, 324(5923):81–85, 2009.
- M. Schumm, B. Eberhard, and P. A. Monkewitz. Self-excited oscillations in the wake of two-dimensional bluff bodies and their control. *J. Fluid Mech.*, 271:17–53, 1994.
- R. Semaan, P. Kumar, M. Burnazzi, G. Tissot, L. Cordier, and B. R. Noack. Reduced-order modelling of the flow around a high-lift configuration with unsteady coanda blowing. *Journal of Fluid Mechanics*, 800:72–110, 8 2016. ISSN 1469-7645. doi: 10.1017/jfm.2016.380. URL [http://journals.cambridge.org/article\\_S0022112016003803](http://journals.cambridge.org/article_S0022112016003803).
- T. K. Sengupta, S. I. Haider, M. K. Parvathi, and G. Pallavi. Enstrophy-based proper orthogonal decomposition for reduced-order modeling of flow past a cylinder. *Physical Review E*, 91(4):043303, 2015.

- D. Sipp and A. Lebedev. Global stability of base and mean flows: a general approach and its applications to cylinder and open cavity flows. *J. Fluid Mech.*, 593:333–358, 2007.
- D. Sipp and P. J. Schmid. Linear closed-loop control of fluid instabilities and noise-induced perturbations: A review of approaches and tools. *Applied Mechanics Reviews*, 68(2):020801, 2016.
- D. Sipp, O. Marquet, P. Meliga, and A. Barbagallo. Dynamics and control of global instabilities in open-flows: a linearized approach. *Applied Mechanics Reviews*, 63(3):030801, 2010.
- L. Sirovich. Turbulence and the dynamics of coherent structures. Part I: Coherent structures. *Quarterly of Applied Mathematics*, 45(3):561–571, 1987.
- G. Tadmor, O. Lehmann, B. R. Noack, and M. Morzyński. Mean field representation of the natural and actuated cylinder wake. *Physics of Fluids (1994-present)*, 22(3):034102, 2010.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.
- J. H. Tu, C. W. Rowley, D. M. Luchtenburg, S. L. Brunton, and J. N. Kutz. On dynamic mode decomposition: theory and applications. *Journal of Computational Dynamics*, 1(2):391–421, 2014.
- W. X. Wang, R. Yang, Y. C. Lai, V. Kovanis, and C. Grebogi. Predicting catastrophes in nonlinear dynamical systems by compressive sensing. *Physical Review Letters*, 106:154101–1–154101–4, 2011.
- S. Wiggins. *Introduction to applied nonlinear dynamical systems and chaos*, volume 2 of *Texts in Applied Mathematics*. Springer Science & Business Media, Berlin, Heidelberg, 2003.
- K. Willcox and J. Peraire. Balanced model reduction via the proper orthogonal decomposition. *AIAA Journal*, 40(11):2323–2330, 2002.
- M. O. Williams, I. G. Kevrekidis, and C. W. Rowley. A data-driven approximation of the Koopman operator: extending dynamic mode decomposition. *Journal of Nonlinear Science*, 2015.
- S. Yamouni, D. Sipp, and L. Jacquin. Interaction between feedback aeroacoustic and acoustic resonance mechanisms in a cavity flow: a global stability analysis. *Journal of Fluid Mechanics*, 717:134–165, 2013.
- A Zebib. Stability of viscous flow past a circular cylinder. *Journal of Engineering Mathematics*, 21(2):155–165, 1987.
- H.-Q. Zhang, U. Fey, B. R. Noack, M. König, and H. Eckelmann. On the transition of the cylinder wake. *Physics of Fluids (1994-present)*, 7(4):779–794, 1995.