

# Constrained sparse Galerkin regression

Jean-Christophe Loiseau<sup>1,†</sup> and Steven 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

(Received 23 December 2016; revised 18 August 2017; accepted 9 November 2017;  
first published online 10 January 2018)

The sparse identification of nonlinear dynamics (SINDy) is a recently proposed data-driven modelling framework that uses sparse regression techniques to identify nonlinear low-order models. With the goal of low-order models of a fluid flow, we combine this approach with dimensionality reduction techniques (e.g. proper orthogonal decomposition) and extend it to enforce physical constraints in the regression, e.g. 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 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 such an approach is demonstrated on two canonical 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 favourably against reduced-order models obtained from a standard Galerkin projection procedure. Finally, the entire code base for our constrained sparse Galerkin regression algorithm is freely available online.

**Key words:** low-dimensional models, nonlinear dynamical systems

---

## 1. Introduction

Fluid flows are characterized by high-dimensional, nonlinear dynamics that gives rise to rich structures. Despite this apparent complexity, the dynamics often evolves 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 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 (Fabbiane *et al.* 2014; Brunton & Noack 2015; Sipp & Schmid 2016; Rowley & Dawson 2017).

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

† Email address for correspondence: [loiseau.jc@gmail.com](mailto:loiseau.jc@gmail.com)

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) (Sirovich 1987; Berkooz, Holmes & Lumley 1993). 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) (Willcox & Peraire 2002; Rowley 2005), have also been widely applied for closed-loop flow control (Ilak & Rowley 2008; Bagheri, Brandt & Henningson 2009; Illingworth, Morgans & Rowley 2010). Although such approaches to model reduction have been widely successful for linear systems, as described in the recent review by Rowley & Dawson (2017) 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 nonlinear dynamics of the von Kármán 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 this reduced-order modelling strategy 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 variety of different techniques such as the eigensystem realization algorithm (ERA) (Juang & Pappa 1985), dynamic mode decomposition (DMD) (Rowley *et al.* 2009; Schmid 2010; Kutz *et al.* 2016), Koopman theory (Mezić 2005, 2013) and variants (Tu *et al.* 2014; Williams, Kevrekidis & Rowley 2015), cluster reduced-order modelling (CROM) (Kaiser *et al.* 2014), NARMAX (Glaz, Liu & Friedmann 2010; Zhang *et al.* 2012; Billings 2013; Semeraro *et al.* 2017) and network analysis of fluids (Nair & Taira 2015). Advances in machine learning are also greatly expanding the ability to extract governing dynamics purely from data. Neural networks have long been used for flow modelling (Milano & Koumoutsakos 2002; Zhang & Duraisamy 2015) and control (Lee *et al.* 1997), and recently deep neural networks (Krizhevsky, Sutskever & Hinton 2012) have been used for Reynolds averaged turbulence modelling (Ling, Kurzawski & Templeton 2016; Kutz 2017). However, many approaches in machine learning, such as neural networks, are prone to overfitting, do not yield interpretable models, and make it difficult to incorporate known physical constraints.

Advanced regression methods from statistics, such as genetic programming or sparse regression, are driving new algorithms that identify parsimonious nonlinear dynamics from measurements of complex systems. Bongard & Lipson (2007) and Schmidt & 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, Proctor & Kutz (2016a) 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 and remaining interpretable.

Many of these regression techniques can be recast into a minimization problem and their solution can be obtained using efficient algorithms available in libraries such as CVXOPT (Andersen, Dahl & Vandenberghe 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 (Noack, Morzynski & Tadmor 2011; Balajewicz, Dowell & Noack 2013; Carlberg, Tuminaro & Boggs 2015; Schlegel & Noack 2015). A notable exception is the physics-constrained multi-level quadratic regression used to identify models in climate and turbulence (Majda & Harlim 2012). Including energy-preserving constraints is known to improve the long-time stability and performance of nonlinear models, while standard Galerkin projection methods often suffer from stability issues (Carlberg, Barone & Antil 2017). Starting from the original SINDy algorithm (Brunton *et al.* 2016a), 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 & Van Loan 2012) to incorporate additional constraints in the SINDy algorithm for the sparse identification of the low-dimensional dynamical system. The ability of the present approach, hereafter named sparse Galerkin regression, is demonstrated on two different flow configurations: the emblematic two-dimensional cylinder flow and the shear-driven cavity flow. We also show that including higher-order nonlinearities in the regression improves the stability and accuracy of resulting models, capturing the effect of truncated low-energy modes on the dynamics of high-energy modes.

The manuscript is organized 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. 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 summarizes our key findings and provides 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 modified version of the sparse identification of nonlinear dynamics (SINDy) method (Brunton *et al.* 2016a). The original SINDy algorithm is introduced in §2.1, and the 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>). 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 from data has been a central challenge in mathematical physics, with a particularly rich history in fluid dynamics. Typically, the structure of the system 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 algorithm (Brunton *et al.* 2016a) bypasses the intractable brute force search through all possible model structures, leveraging the observation that many dynamical systems

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) \tag{2.1}$$

have dynamics  $\mathbf{f}$  that is 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. Note that the vector  $\mathbf{x}$  refers to the state of the system, and may be replaced with a vector of POD coefficients  $\mathbf{a}$  for reduced-order modelling.

First, time-series data are collected from (2.1) and formed into a data matrix,

$$\mathbf{X} = [\mathbf{x}(t_1) \quad \mathbf{x}(t_2) \quad \cdots \quad \mathbf{x}(t_m)]^T, \tag{2.2}$$

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

$$\dot{\mathbf{X}} = [\dot{\mathbf{x}}(t_1) \quad \dot{\mathbf{x}}(t_2) \quad \cdots \quad \dot{\mathbf{x}}(t_m)]^T. \tag{2.3}$$

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:

$$\Theta(\mathbf{X}) = [\mathbf{1} \quad \mathbf{X} \quad \mathbf{X}^2 \quad \cdots \quad \mathbf{X}^d \quad \cdots \quad \sin(\mathbf{X}) \quad \cdots]. \tag{2.4}$$

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 (2.1) may now be represented in terms of the data matrices in (2.3) and (2.4) as

$$\dot{\mathbf{X}} = \Theta(\mathbf{X}). \tag{2.5}$$

Each column  $\boldsymbol{\Xi}_k$  in  $\boldsymbol{\Xi}$  is a vector of coefficients determining the active terms in the  $k$ th row equation in (2.1). A parsimonious model will provide an accurate model fit in (2.5) with as few terms as possible in  $\boldsymbol{\Xi}$ . Such a model may be identified using a convex  $\ell_1$ -regularized sparse regression:

$$\boldsymbol{\Xi}_k = \underset{\boldsymbol{\Xi}'_k}{\operatorname{argmin}} \|\dot{\mathbf{X}}_k - \Theta(\mathbf{X})\boldsymbol{\Xi}'_k\|_2 + \lambda \|\boldsymbol{\Xi}'_k\|_1. \tag{2.6}$$

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). Alternatively, symbolic regression techniques such as the fast function extraction could be used to identify nonlinear terms in the dynamics (McConaghy 2011). Note that the reformulation of the nonlinear dynamics into a linear regression framework via a library of candidate basis functions in (2.5) was

developed earlier by Yao & Bollt (2007), although they did not obtain parsimonious models with sparsity promoting techniques.

The sparse vectors  $\mathbf{E}_k$  may be synthesized into a nonlinear dynamical system model:

$$\dot{x}_k = \Theta(\mathbf{x}) \mathbf{E}_k. \quad (2.7)$$

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. Note that, if  $\Theta(\mathbf{X})$  consists only of linear terms, and if the sparsity promoting term is set to  $\lambda = 0$ , this algorithm reduces to dynamic mode decomposition (Rowley *et al.* 2009; Schmid 2010; Kutz *et al.* 2016). A major benefit of the SINDy architecture is its ability to identify parsimonious models that contain only the required nonlinear terms, resulting in interpretable models and avoiding overfitting.

Recent extensions 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.* 2017; Schaeffer 2017), and has been extended to include inputs and control (Brunton, Proctor & Kutz 2016b). Other nonlinear modelling techniques, such as NARMAX (Billings 2013), have been widely applied to problems in fluid mechanics, including modelling of wave packets in a turbulent jet (Semeraro *et al.* 2017), flutter instability (Zhang *et al.* 2012) and unsteady aerodynamics (Glaz *et al.* 2010). SINDy is closely related to NARMAX modelling, which has also been extended to include sparsity-promoting techniques such as the LASSO for parsimonious modelling (Kukreja, Löfberg & Brenner 2006; Kukreja & Brenner 2007; Linscott & Wiklund 2014). However, the extensions of SINDy to identify partial differential equations (Rudy *et al.* 2017; Schaeffer 2017) and biological regulatory networks (Mangan *et al.* 2016) highlight the flexibility of this simple regression framework.

## 2.2. Constrained sparse identification

It has been shown in § 2.1 that the identification problem can be cast as a convex optimization problem where the sparsity of the solution  $\mathbf{E}$  can be promoted using an  $\ell_1$  penalization. Alternatively, sparsity can also be promoted by using the sequential thresholded least-squares algorithm as in Brunton *et al.* (2016a). In this case, the convex minimization problem can be rewritten as

$$\left. \begin{aligned} \min_{\mathbf{E}} \|\Theta(\mathbf{X})\mathbf{E} - \dot{\mathbf{X}}\|_2^2, \\ \text{subject to } \mathbf{C}\boldsymbol{\xi} = \mathbf{d}, \end{aligned} \right\} \quad (2.8)$$

where  $\boldsymbol{\xi} = \mathbf{E}(\cdot)$  is the vectorized form of the sparse matrix of coefficients, and where  $\mathbf{C}\boldsymbol{\xi} = \mathbf{d}$  are linear equality constraints, which can be used to enforce that some entries of  $\boldsymbol{\xi}$  are equal to zero. The minimization 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.

From a practical point of view, each iteration of (2.8) can be recast as an unconstrained problem using an augmented functional formulation where the constraints are imposed via Lagrange multipliers. The resulting unconstrained minimization problem then reads

$$\min_{\xi, z} \|\Theta(\mathbf{X})\boldsymbol{\Xi} - \dot{\mathbf{X}}\|_2^2 + z^T(\mathbf{C}\xi - d). \tag{2.9}$$

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

$$\begin{bmatrix} 2\hat{\Theta}(\mathbf{X})^T \hat{\Theta}(\mathbf{X}) & \mathbf{C}^T \\ \mathbf{C} & 0 \end{bmatrix} \begin{bmatrix} \xi \\ z \end{bmatrix} = \begin{bmatrix} 2\hat{\Theta}(\mathbf{X})^T \dot{\mathbf{X}}(\cdot) \\ d \end{bmatrix}, \tag{2.10}$$

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  $\boldsymbol{\Xi}$  into  $\xi = \boldsymbol{\Xi}(\cdot)$ ) 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  $[\hat{\Theta}(\mathbf{X}) \ \mathbf{C}]^T$  has full column rank.

Interestingly, the linear equality constraints  $\mathbf{C}\xi = d$  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. Specific constraints required to conserve energy in a fluid are derived later in §3.

*Notes on the numerical implementation.* Although it has been extended with the possibility of including user-provided constraints, SINDy is at its core a classical linear regression problem. Its computational cost depends essentially on:

- (i) the dimension of the state vector  $\mathbf{x}$  characterizing the system,
- (ii) the number of functions included in the pool of admissible functions  $\Theta(\mathbf{x})$ ,
- (iii) the algorithm used to solve the optimization problem.

Since the systems considered in the present work are characterized by only three degrees of freedom and only have up to 20 different functions included in  $\Theta(\mathbf{x})$ , the solution to the optimization problem has been obtained directly using the closed-form solution of the KKT equations, which involves the inversion of a  $60 \times 60$  symmetric-positive-definite matrix. If the number of degrees of freedom and/or the number of admissible functions considered is relatively large, the constrained optimization problem can be solved using gradient descent algorithm or a variant, e.g. L-BFGS or stochastic gradient descent. In this work, the constrained sparse regression algorithm is implemented in Python. It uses the CVXOPT library (Andersen *et al.* 2013) to solve the constrained least-squares problem. Moreover, every time



an additional sparsity constraint is added as a new row to the matrix  $\mathbf{C}$ , a QR rank-revealing decomposition of  $\mathbf{C}$  is performed using SciPy (Jones *et al.* 2001) to ensure it has full rank (i.e. no linearly dependent constraints).

### 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

$$\frac{\partial \mathbf{u}}{\partial t} = -(\mathbf{U}_b \cdot \nabla) \mathbf{u} - (\mathbf{u} \cdot \nabla) \mathbf{U}_b - \nabla p + \frac{1}{Re} \nabla^2 \mathbf{u} - (\mathbf{u} \cdot \nabla) \mathbf{u}, \quad (3.1)$$

$$\nabla \cdot \mathbf{u} = 0, \quad (3.2)$$

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 modelling is to obtain a low-dimensional system of the form

$$\frac{d\mathbf{a}}{dt} = \mathcal{L}\mathbf{a} + \mathcal{N}(\mathbf{a}), \quad (3.3)$$

where  $\mathbf{a}$  is a vector of POD coefficients that represent the degrees of freedom of the reduced-order model, and where  $\mathcal{L}$  and  $\mathcal{N}(\mathbf{a})$  are low-dimensional approximation of the linearized Navier–Stokes operator and of the quadratic nonlinear term, respectively.

For the reduced-order model (3.3) to be a good approximation of its high-dimensional counterpart, the former needs to have the same physical properties as the latter. While this may be enforced when the reduced-order model is derived based on a Galerkin projection (Noack *et al.* 2011; Balajewicz *et al.* 2013; Carlberg *et al.* 2015; Schlegel & Noack 2015), these properties need to be actively enforced when a system identification approach such as SINDy is used. This discussion on the different constraints used in the present work rests on the assumption that the library of candidate functions used in the identification is given by

$$\Theta(\mathbf{a}) = [P_1(\mathbf{a}) \quad P_2(\mathbf{a}) \quad \cdots \quad P_N(\mathbf{a})], \quad (3.4)$$

where  $P_i(\mathbf{a})$  defines all the polynomials of degree  $i$  in the entries of  $\mathbf{a}$ . Thus, SINDy models will be obtained in terms of the vector  $\mathbf{a}$  of POD coefficients.

#### 3.1. Constraining the quadratic nonlinear term

The nonlinear Navier–Stokes equations (3.2) are partial differential equations characterized by the quadratic nonlinear term  $-(\mathbf{u} \cdot \nabla) \mathbf{u}$ . It can be shown that

$$\int_{\Omega} \mathbf{u} \cdot (\mathbf{u} \cdot \nabla) \mathbf{u} \, d\Omega = 0, \quad (3.5)$$

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 and parsimony. 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 length scales 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

$$\mathbf{a} \cdot \mathcal{N}(\mathbf{a}) = 0. \tag{3.6}$$

In the rest of this work, all of the identified models will be characterized by three degrees of freedom so that the state vector is given by

$$\mathbf{a} = [a_1 \quad a_2 \quad a_3]^T. \tag{3.7}$$

Expanding (3.6) in terms of the regression coefficients  $\xi$  yields

$$0 = [a_1 \quad a_2 \quad a_3] \begin{bmatrix} \xi_4^{(a_1)} a_1 & \xi_5^{(a_1)} a_1 + \xi_7^{(a_1)} a_2 & \xi_6^{(a_1)} a_1 + \xi_9^{(a_1)} a_3 \\ \xi_4^{(a_2)} a_1 + \xi_5^{(a_2)} a_2 & \xi_7^{(a_2)} a_2 & \xi_8^{(a_2)} a_2 + \xi_9^{(a_2)} a_3 \\ \xi_4^{(a_3)} a_1 + \xi_6^{(a_3)} a_3 & \xi_7^{(a_3)} a_2 + \xi_8^{(a_3)} a_3 & \xi_9^{(a_3)} a_3 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} \\ + [a_1 \quad a_2 \quad a_3] \begin{bmatrix} \xi_8^{(a_1)} a_2 a_3 \\ \xi_6^{(a_2)} a_1 a_3 \\ \xi_5^{(a_3)} a_1 a_2 \end{bmatrix}. \tag{3.8}$$

For (3.8) to hold, the matrix involved in the first term is required to be skew symmetric, while the second term implies  $\xi_8^{(a_1)} + \xi_6^{(a_2)} + \xi_5^{(a_3)} = 0$ . Overall, this gives rise to the following ten different linear equality constraints:

$$\left. \begin{aligned} \xi_8^{(a_1)} + \xi_6^{(a_2)} + \xi_5^{(a_3)} &= 0, \\ \xi_4^{(a_1)} = \xi_7^{(a_2)} = \xi_9^{(a_3)} &= 0, \\ \xi_5^{(a_1)} &= -\xi_4^{(a_2)}, \\ \xi_7^{(a_1)} &= -\xi_5^{(a_2)}, \\ \xi_6^{(a_1)} &= -\xi_4^{(a_3)}, \\ \xi_9^{(a_1)} &= -\xi_6^{(a_3)}, \\ \xi_8^{(a_2)} &= -\xi_7^{(a_3)}, \\ \xi_9^{(a_2)} &= -\xi_8^{(a_3)}, \end{aligned} \right\} \tag{3.9}$$

which induce a coupling of the different ordinary differential equations governing the evolution of  $a_1$ ,  $a_2$  and  $a_3$ . If the system we aim to identify has more degrees of freedom, the exact same procedure applies, although it would require more calculations to derive all of the required constraints.

### 3.2. Including higher-order nonlinearities

Reduced-order modelling based on Galerkin projection usually requires a relatively large projection basis. Despite the low-dimensional effective dynamics 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 reasonably 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. In order to ease the rest of the discussion, let us consider the following generalized mean-field model:

$$\left. \begin{aligned} \frac{da_1}{dt} &= \sigma a_1 + \omega a_2 - a_1 a_3, \\ \frac{da_2}{dt} &= -\omega a_1 + \sigma a_2 - a_2 a_3, \\ \frac{da_3}{dt} &= \lambda(-a_3 + a_1^2 + a_2^2), \end{aligned} \right\} \quad (3.10)$$

with  $\sigma$  and  $\lambda$  being positive constants. Assuming  $\lambda \gg \sigma$  implies that the dynamics of  $a_3$  is entirely enslaved to that of  $a_1$  and  $a_2$ . Here,  $a_3$  thus plays the role of the higher-order POD modes. Using adiabatic elimination (Haken 1983) or centre manifold reduction (Wiggins 2003; Carini, Auteri & Giannetti 2015), it is well known that these enslaved degrees of freedom can be reduced out of the problem, while their influence onto the driving modes can be accounted for by appropriately modifying the nonlinear terms. In the present case,  $a_3(t)$  can be approximated as

$$a_3(t) \approx a_1^2 + a_2^2. \quad (3.11)$$

Inserting this approximation into our original system (3.10), one can recast it as an effective two-dimensional dynamical system given by

$$\left. \begin{aligned} \frac{da_1}{dt} &= \sigma a_1 + \omega a_2 - (a_1^2 + a_2^2)a_1, \\ \frac{da_2}{dt} &= -\omega a_1 + \sigma a_2 - (a_1^2 + a_2^2)a_2. \end{aligned} \right\} \quad (3.12)$$

As can be seen, the influence of the eliminated degree of freedom is accounted for by transforming the original quadratic nonlinearity into an effective cubic one. The same approach has been 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. Such an approach, which can be summarized as derive then reduce, 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 only 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 effectiveness of Galerkin regression, 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 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 (Zebib 1987; Schumm, Eberhard & Monkewitz 1994) and below the onset of three-dimensional instabilities (Zhang *et al.* 1995; Barkley & Henderson 1996). The saturation process of the instability is well described by the first-order self-consistent model of Mantič-Lugo, Arratia & Gallaire (2014). In the fluid dynamics community, a large body of literature exists in which this particular set-up has been chosen to illustrate modal decomposition (Bagheri 2013; Noack *et al.* 2016) and model identification techniques (Noack *et al.* 2003; Sengupta *et al.* 2015; Brunton *et al.* 2016a; Rowley & Dawson 2017). This set-up 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 is governed by the incompressible Navier–Stokes equations. The centre of the cylinder has been chosen as the origin of the reference frame  $(x, y)$ , where  $x$  denotes the streamwise coordinate and  $y$  denotes the spanwise coordinate. This study considers the same computational domain as in Noack *et al.* (2003), extending from  $x = -5$  to  $x = 15$  in the streamwise direction, and from  $y = -5$  to  $y = 5$  in the spanwise direction. A uniform velocity profile is prescribed at the inflow, a classical stress-free boundary condition is used at the outflow, and free-slip boundary conditions are used on the lateral boundaries of the computational domain. Based on the spectral element solver Nek5000 (Fischer, Lottes & Kerkemeir 2008), the domain is discretized by 1832 seventh-order spectral elements. Finally, the time integration of the diffusive terms relies on a backward differentiation of order 3, while the convective terms are advanced in time based on a third-order accurate extrapolation.

The vorticity field of the linearly unstable fixed point  $U_b$ , computed using the selective frequency damping approach (Åkervik *et al.* 2006), is shown in figure 1(b). Figure 1(a,c) also provides the eigenspectrum of the linearized Navier–Stokes operator and the vorticity field associated with 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 nonlinearly saturated von Kármán vortex street (Barkley 2006).

In the rest of this work, three different transient evolutions are considered. The first one, shown in figure 3(a), is started with the initial condition

$$U(\mathbf{x}, 0) = U_b(\mathbf{x}) + \epsilon \mathbf{u}'(\mathbf{x}), \quad (4.1)$$

where  $U_b$  is the linearly unstable base flow,  $\mathbf{u}'$  is the leading unstable eigenmode normalized such that it is unit norm and  $\epsilon = 10^{-6}$ . A direct numerical simulation has been run until a statistical steady state has been achieved. The dynamics of the system on the final attractor is then equidistantly sampled using  $M = 1000$  velocity field snapshots with a sampling frequency approximately 30 times larger than the vortex shedding frequency (Noack *et al.* 2016). The shift mode, denoted  $\mathbf{u}_\Delta$  and depicted in figure 2(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 modelling (Noack *et al.* 2003; Tadmor *et al.* 2010). The snapshot POD method

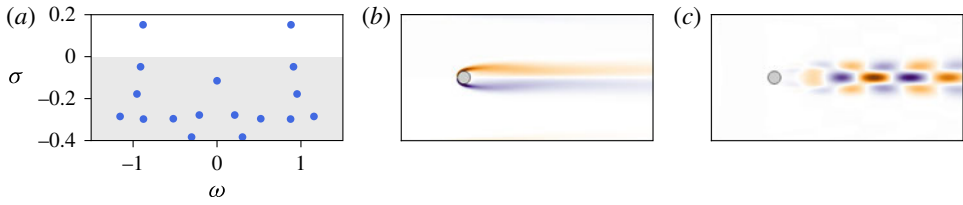


FIGURE 1. (Colour online) (a) Eigenspectrum of the linearized Navier–Stokes operator for the two-dimensional cylinder flow at  $Re = 100$ . Vorticity fields of (b) the base flow and (c) the leading linearly unstable eigenmode.

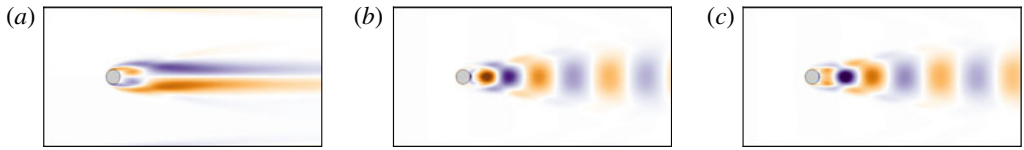


FIGURE 2. (Colour online) Vorticity fields of (a) the shift mode, (b) the first and (c) second POD modes of the cylinder flow at  $Re = 100$ .

of Sirovich (1987) has then been used to extract the two most energetic modes  $\mathbf{u}_1$  and  $\mathbf{u}_2$ , depicted in figures 2(b) and 2(c), respectively. The evolution in time of the POD coefficients is shown in figure 3(a), along with a projection of the system's trajectory onto the  $(a_1, a_\Delta)$  plane, where  $a_1(t)$  is the amplitude of the POD mode  $\mathbf{u}_1$  and  $a_\Delta(t)$  is the amplitude of the shift mode  $\mathbf{u}_\Delta$ . As shown in Noack *et al.* (2003), the system evolves on a low-dimensional paraboloid manifold characterized by  $a_\Delta \propto a_1^2 + a_2^2$ .

Two additional transient evolutions, started with the initial conditions

$$\mathbf{U}(\mathbf{x}, t) = \mathbf{U}_b \pm 2.25\mathbf{u}_\Delta, \quad (4.2)$$

are considered in order to capture the off-manifold dynamics. The first one, shown in figure 3(b), corresponds to a direct numerical simulation started from the mean flow. The second one, shown in figure 3(c), is characterized by an initial condition having a reversed flow region longer than that of the linearly unstable base flow. In both cases, the flow is rapidly attracted toward the vicinity of the fixed point before escaping away from it due to its linearly unstable nature. Including this off-manifold dynamics was deemed necessary in order to identify a physically consistent equation governing the dynamics of the mean-flow distortion  $a_\Delta(t)$ . In the rest, the transient evolutions shown in figure 3(a,b) are forming the training dataset used for the identification. The remaining transient evolution (see figure 3c) is used for cross-validation purposes.

#### 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 localized 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 perturbations to grow as they are convected downstream. On the other hand, the inner-cavity

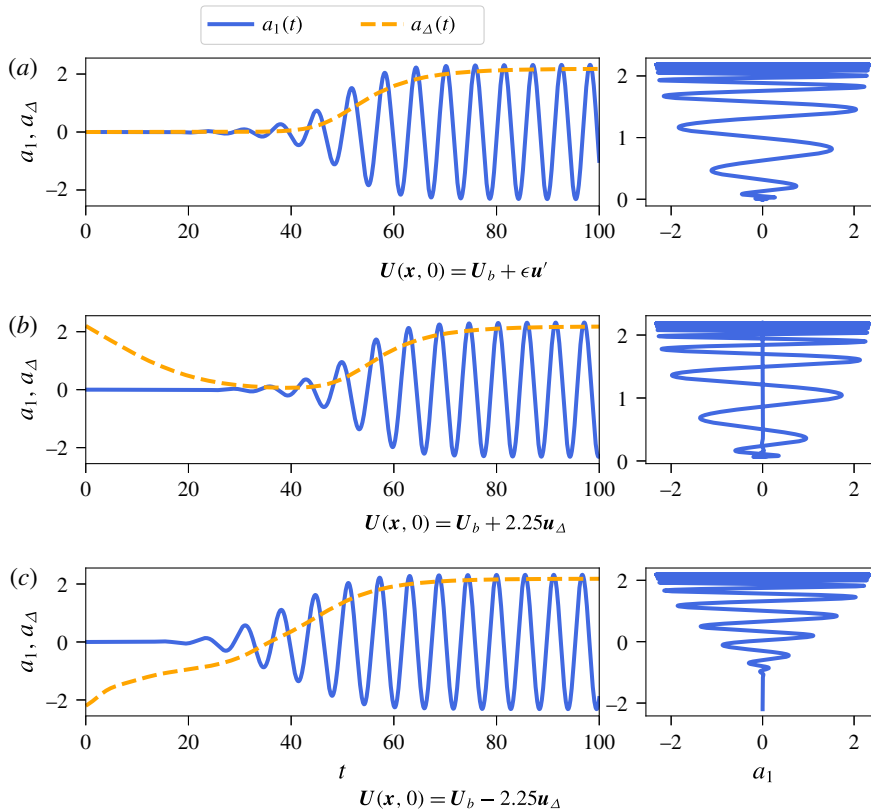


FIGURE 3. (Colour online) Time evolution of the POD coefficients and corresponding trajectory in the phase space projected onto the  $(a_1, a_\Delta)$  plane for the cylinder flow at  $Re = 100$ . The time evolution of  $a_2(t)$ , not shown, is very similar to that of  $a_1(t)$ . All three datasets used in the present work are presented.

recirculating flow and the instantaneous pressure feedback provide the mechanisms allowing these same perturbations to eventually re-excite the upstream shear layer. The coupling between these mechanisms gives rise to a linearly unstable feedback loop at sufficiently high Reynolds numbers. In this case, the overall saturation process of the instability is well described by the second-order self-consistent model recently proposed by Meliga (2017). 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, Colonius & Basu 2002; Yamouni, Sipp & Jacquin 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, Sipp & Schmid 2009), investigation of the nonlinear saturation process of globally unstable flows (Sipp & Lebedev 2007; Meliga 2017) 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 & 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 4 for the sake of completeness.

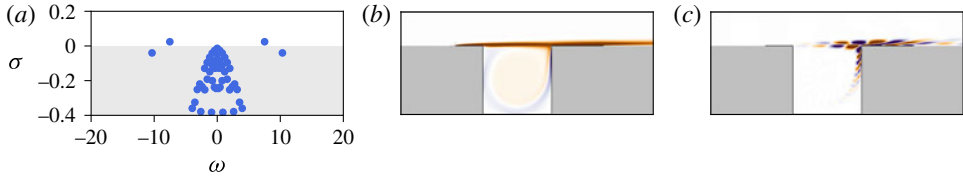


FIGURE 4. (Colour online) (a) Eigenspectrum of the linearized Navier–Stokes operator for the shear-driven cavity flow at  $Re = 4250$ . Vorticity fields of (b) the base flow and (c) the leading linearly unstable eigenmode. The dashed lines indicate the spatial extent over which the free-slip boundary condition is imposed.

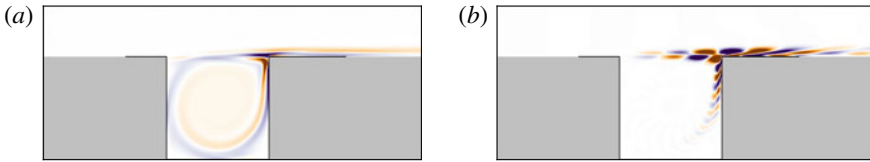


FIGURE 5. (Colour online) Vorticity fields of (a) the shift mode and (b) the first POD mode for the shear-driven cavity flow at  $Re = 4250$ .

Three different transients are once again considered. The first one, shown in figure 6(a), is started with the initial condition

$$\mathbf{U}(\mathbf{x}, 0) = \mathbf{U}_b(\mathbf{x}) + \epsilon \mathbf{u}'(\mathbf{x}), \quad (4.3)$$

where  $\mathbf{U}_b$  is the linearly unstable base flow,  $\mathbf{u}'$  is the leading unstable eigenmode normalized such that it is unit norm and  $\epsilon = 10^{-8}$ . This direct numerical simulation has been run until a statistically steady state has been reached. As for the cylinder flow, the dynamics of the attractor has been equidistantly sampled using  $M = 1000$  velocity field snapshots with a sampling frequency approximately 30 times larger than the oscillation frequency of the shear layer. The shift mode  $\mathbf{u}_\Delta$  is depicted in figure 5(a) and the leading POD mode is shown in figure 5(b). 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  $Re = 4250$ . Note furthermore that, 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 figures 3(b) and 6(b). As for the cylinder, two additional transients, started from either side of the fixed point in the direction of the shift mode, have been included. Once again, the transients shown in figure 6(a,b) are forming the training dataset used for the identification problem, while the transient depicted in figure 6(c) will be used for cross-validation purposes only.

## 5. Results and discussion

Following the seminal work of Noack *et al.* (2003), so-called quadratic Galerkin regression models are constructed from the basic building blocks necessary for reduced-order modelling of the flow configurations considered, i.e. a linear operator

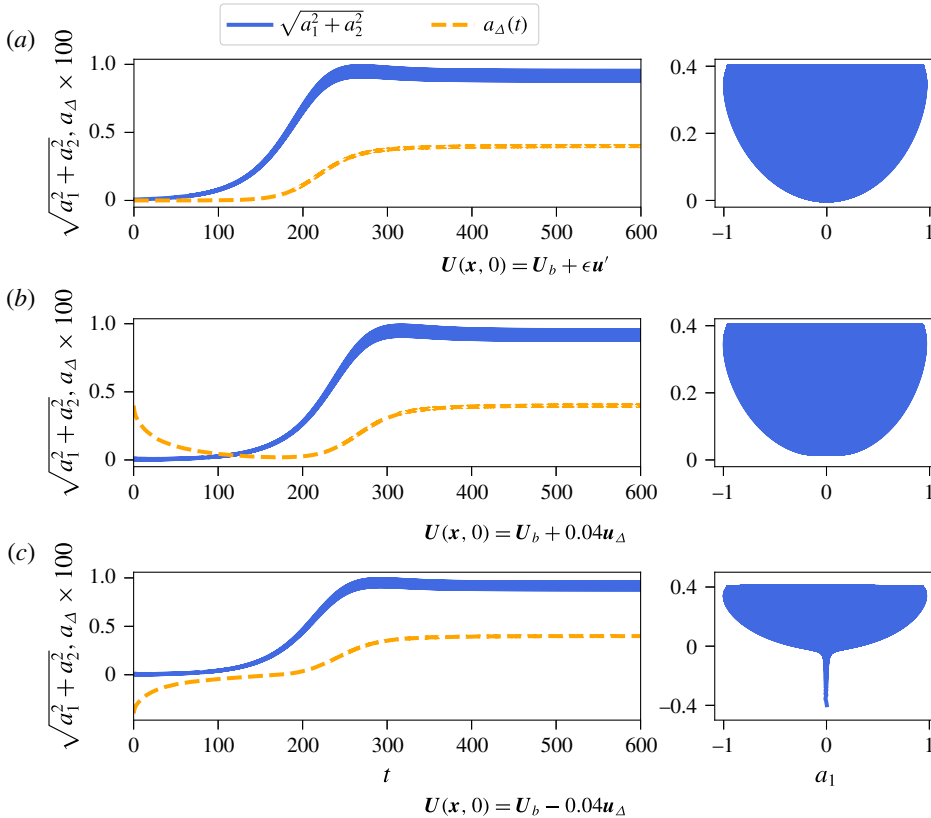


FIGURE 6. (Colour online) Time evolution of the POD coefficients and corresponding trajectory in the phase space projected onto the  $(a_1, a_2)$  plane for the shear-driven cavity flow at  $Re = 4250$ . The time evolution of  $a_2(t)$ , not shown, is very similar to that of  $a_1(t)$ . All three datasets used in the present work are presented.

$\mathcal{L}$  and an energy-preserving quadratic nonlinearity  $\mathcal{N}(\mathbf{a})$ . For that purpose, the library  $\Theta(\mathbf{a})$  used in the identification process is defined as  $[P_1(\mathbf{a}) \ P_2(\mathbf{a})]$ , i.e. all the polynomials of degree 2 or less in the entries of  $\mathbf{a}$ . 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  $[P_1(\mathbf{a}) \ P_2(\mathbf{a}) \ P_3(\mathbf{a})]$ , 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.

### 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:

- (i) the minimal Galerkin projection model including only the shift mode and the first two POD modes,

- (ii) a Galerkin projection model including the shift mode and the first eight POD modes.

### 5.1.1. Model selection

Model selection and cross-validation are crucial components of system identification. The goal is to identify, among all candidate models, the parsimonious model that optimally balances model accuracy and model complexity. As the sparsifying parameter  $\lambda$  is varied in the SINDy procedure, a Pareto front is swept out, reducing the combinatorially many candidate models down to a small handful of candidate models. Mangan *et al.* (2017) have recently demonstrated how SINDy can be combined with the well-known Akaike information criterion (AIC) (Akaike 1974) or the Bayes information criterion (BIC) (Schwarz *et al.* 1978) in order to select the most parsimonious model from this Pareto front. Given a candidate model, the associated AIC score is given by

$$\text{AIC} = 2k - 2 \ln(L(\mathbf{a}, \boldsymbol{\Xi})) + 2 \frac{(k+1)(k+2)}{(m-k-2)}, \quad (5.1)$$

where  $L(\mathbf{a}, \boldsymbol{\Xi})$  is the loss function of the observations  $\mathbf{a}$  given the best-fit parameter values  $\boldsymbol{\Xi}$  of the candidate model and  $k$  is the total number of free parameters. The last term in (5.1) is a finite sample size correction where  $m$  is the total number of observations used to cross-validate the model. These training and testing datasets are the same as those shown in figure 3. For two models of the same accuracy, the AIC score will penalize the one having the larger number of terms.

The AIC scores for each candidate model can have a wide range of values, hence requiring a rescaling by the minimum AIC value. The relative AIC score is thus given by

$$\Delta = \text{AIC} - \text{AIC}_{\min}. \quad (5.2)$$

The different candidate models can then be ranked based on this relative AIC score. Following Mangan *et al.* (2017), models with  $\Delta \leq 2$  have so-called strong support, models with  $4 \leq \Delta \leq 7$  have weak support and models with  $\Delta \geq 10$  have no support. It should be emphasized that the model characterized by  $\Delta = 0$  is not necessarily the best model possible, but only the best one among the different models tested.

Figure 7 depicts the distribution of all the different models identified in the complexity versus  $\text{AIC}_c$  plane for the two-dimensional cylinder flow at  $Re = 100$ . Note that none of the unconstrained models are shown as they have all diverged in the cross-validation stage when trying to reproduce the dynamics of the transient evolution shown in figure 3(c). This observation will be investigated in § 5.1.3. Although the quadratic Galerkin regression models have lower complexity and appear at first to be more physical, it is interesting to note that they are largely dominated by the cubic models. In the rest of this section, only the best constrained quadratic model and the best constrained cubic model are considered. Excluding the transient evolution shown in figure 3(c) from the validation stage, a similar analysis is performed in order to select the unconstrained models considered in the following sections.

### 5.1.2. Qualitative comparisons

Let us consider the first transient evolution forming our training dataset, shown in figure 3(a). Figures 8 and 9 provide visual comparisons of the dynamics predicted by the different models. The time evolution of the mean-flow distortion predicted by the low-dimensional Galerkin projected systems, shown in figure 8(a), indicates that the duration of the transients is largely over-estimated and that an energy overshoot



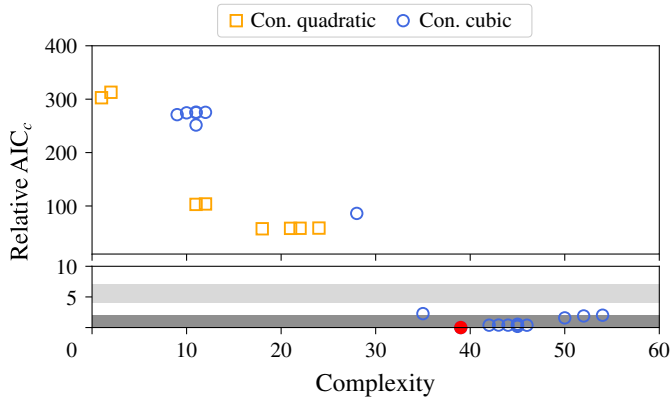


FIGURE 7. (Colour online) Distribution of all the different constrained low-order models identified in the complexity versus  $AIC_c$  plane. The red circle corresponds to the best model identified.

occurs once nonlinear saturation kicks in. On the one hand, the over-estimation of the duration of transients results from the fact that the leading POD modes (see figure 2) provide only a crude approximation of the leading linear instability eigenmodes (see figure 1). On the other hand, the 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 entirely 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 constrained and unconstrained quadratic Galerkin regression models suffer 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. Moreover, the unconstrained model *a priori* performs better than the constrained one. Comparatively, the cubic Galerkin regression models provide an almost perfect fit to the original data, as shown in figures 8(b) and 9(d): 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%. Moreover, the inclusion of the cubic nonlinearities, modelling the influence of the truncated modes, has a stabilizing effect, hence preventing the energy overshoot and larger limit cycle amplitude observed for the quadratic models.

### 5.1.3. Quantitative analysis

Table 1 reports the growth rate and circular frequency of the leading eigenvalue of the low-dimensional linear operator  $\mathcal{L}$  for the different models considered and compares it against the values obtained from a linear stability analysis of the linearized Navier–Stokes equations. As discussed previously, the leading POD modes (see figure 2) provide only a crude approximation of the leading linear instability eigenmodes (see figure 1). As a result of this crude approximation, the growth rate of the leading eigenvalue of the low-dimensional operator  $\mathcal{L}$  obtained by Galerkin projection is three to four times smaller than the actual value obtained by linear stability analysis, hence explaining the over-estimation of the transients duration, while the associated frequency is 10% larger than the actual one. Similarly, the

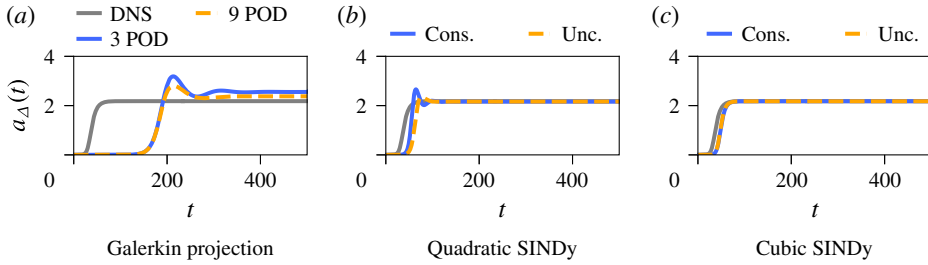


FIGURE 8. (Colour online) Comparison of the time evolution of the mean-flow distortion  $a_\Delta$  predicted by the different data-driven models for the two-dimensional cylinder flow at  $Re = 100$ . The transient evolution considered, which is part of the training dataset, is the one depicted in figure 3(a).

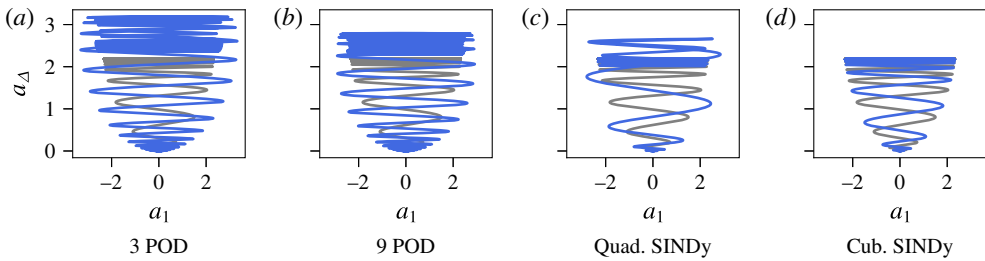


FIGURE 9. (Colour online) Comparison of the trajectory in the  $a_1$ - $a_\Delta$  plane predicted by the different reduced-order models for the two-dimensional cylinder flow at  $Re = 100$ . The light grey trajectory is the one given by a direct numerical simulation. Only the constrained cases are reported for the Galerkin regression models.

growth rate of the unconstrained and constrained quadratic models are 30 % smaller than the correct value. Compared to these models, the spectral properties of the cubic Galerkin regression models are in excellent agreement with the results obtained from linear stability analysis of the Navier–Stokes equations, the identified growth rate  $\sigma$  being only 6 % larger than the true one while the corresponding frequency is only up to 1 % different than the correct value. Introducing higher-order nonlinearities in the model identification thus not only allows us to take into account the influence of the truncated modes on the driving ones, but it also enables the optimization procedure to correctly estimate the growth rate of the linear instability.

Now focusing our attention on the decay rate  $\sigma_\Delta$  of the shift mode, it can be seen that both constrained models identify the mean-flow distortion as being linearly stable, so does the unconstrained cubic model. Although the unconstrained quadratic model appears to outperform the constrained one when looking at the evolution depicted in figure 8(b), it can be seen from table 1 that it surprisingly identifies the mean-flow distortion as being linearly unstable. This misprediction of the linearly stable nature of the mean-flow distortion also explains why the present unconstrained model fails to reproduce the dynamics of the third transient evolution (see figure 3c) used in the cross-validation stage. If one were to consider only our first two transient evolutions (figure 3a,b) without prior knowledge of the problem, one could easily conclude that the system is indeed 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

	Galerkin projection	Unc. quad. SINDy	Con. quad. SINDy	Unc. cubic SINDy	Con. cubic SINDy	N–S
$\sigma$	0.044	0.111	0.115	0.156	0.162	0.152
$\omega$	0.97	0.81	0.81	0.87	0.88	0.88
$\sigma_{\Delta}$	-0.04	0.032	-0.122	-0.052	-0.122	N/A

TABLE 1. Comparison of growth rate  $\sigma$  and circular frequency  $\omega$  of the leading eigenvalue for the different models considered. Results from a global stability analysis of the Navier–Stokes (N–S) equations are also reported. The last row also reports the value of the decay rate  $\sigma_{\Delta}$  associated with the shift mode. Note that only the first two transients shown in figure 3 are part of the training dataset.

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 measurements available in our training dataset. Such a model is however not acceptable as it could lead to a misunderstanding of the physics at play. This example thus clearly demonstrates the benefits of introducing physics into the identification process: 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.

Finally, figure 10 depicts time series of  $a_1(t)$  in the nonlinearly saturated stage and the associated Fourier spectrum for the different models considered. Comparing these different Fourier spectra, it is clear that the vortex shedding frequency  $\omega = 1.15$  ( $St = 0.18$ ) predicted by all the models in the nonlinear regime is in excellent agreement with that observed from direct numerical simulation.

### 5.2. Shear-driven cavity flow

Figure 11 provides 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 shear-driven cavity flow at  $Re = 4250$ . It also provides the dynamics predicted by two additional data-driven reduced-order models, namely:

- (i) the minimal Galerkin projection model including only the shift mode and the first two POD modes,
- (ii) a Galerkin projection model including the shift mode and the first six POD modes.

The model selection procedure, being the same as described in § 5.1, is thus not discussed again for the present case.

#### 5.2.1. Overview

Figure 11 provides visual comparisons of the dynamics predicted by the different models for the two-dimensional shear-driven cavity flow at  $Re = 4250$ . Although the geometry and the physics are quite different from that of the two-dimensional cylinder flow, the present Galerkin projection models suffer from similar drawbacks as before: a misprediction of the transients duration and the saturation to higher

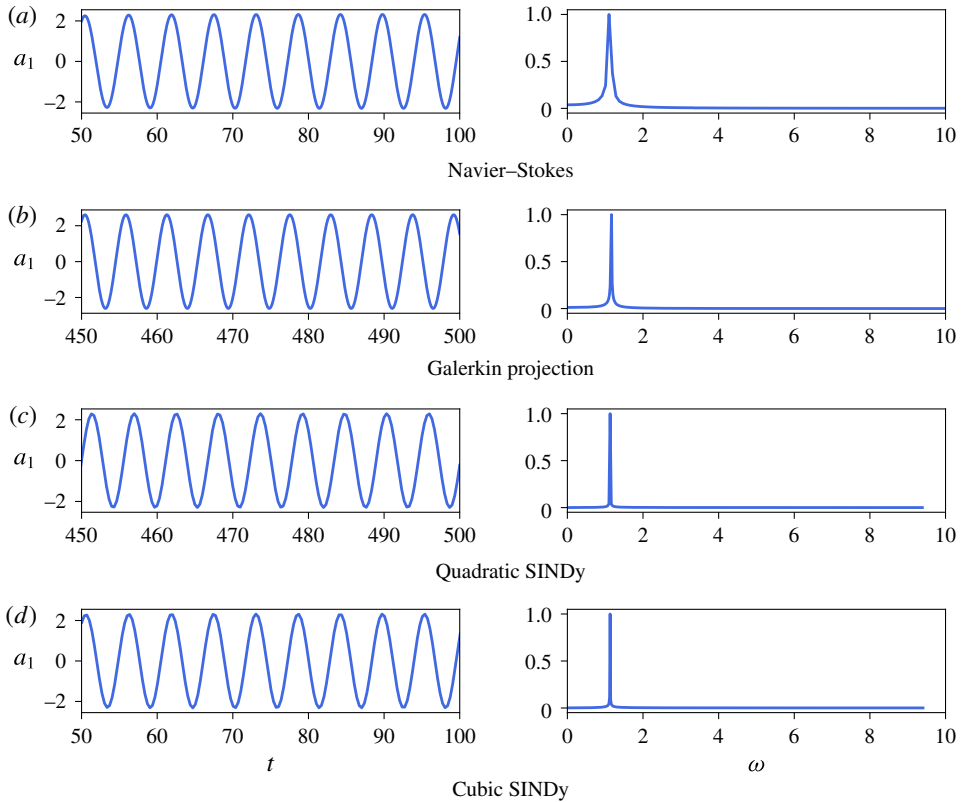


FIGURE 10. (Colour online) Time series of the nonlinearly saturated dynamics of  $a_1(t)$  and normalized Fourier spectra for the different models considered for the two-dimensional cylinder flow at  $Re = 100$ . Note that the larger width of the peak in (a) is solely related to the smaller integration window used in direct numerical simulation. Only the constrained cases are reported for the Galerkin regression models.

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. Looking now at figure 11(b), the two quadratic Galerkin regression models correctly reproduce the asymptotic dynamics of the shear-driven cavity flow. Both of them slightly over-predict the duration of the transients. Finally, both cubic models appear to exhibit similar accuracy as shown in figure 11(c), although the unconstrained version appears to saturate slightly faster.

### 5.2.2. Quantitative analysis

Table 2 provides a comparison of the growth rate  $\sigma$  and circular frequency  $\omega$  of the leading unstable eigenvalue for each of the models considered. As assessed from figure 11, all of these growth rates are in qualitative agreement with that obtained from a global linear stability analysis of the Navier–Stokes equations, the constrained cubic model differing by less than 1%. One way to further improve the accuracy of the quadratic models would be to constrain the eigenspectrum of the low-dimensional linear operator to be a subset of its high-dimensional counterpart. Such a constraint

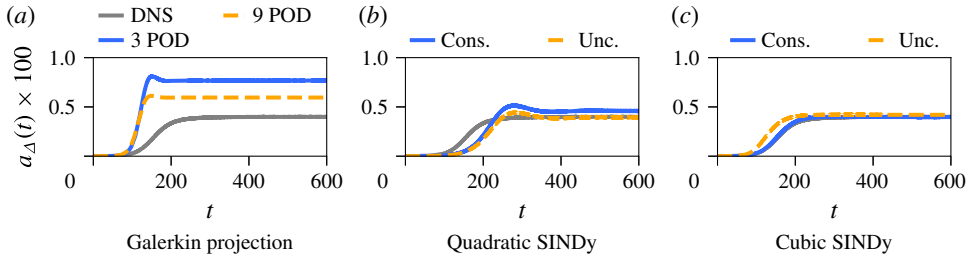


FIGURE 11. (Colour online) Comparison of the time evolution of the mean-flow distortion  $a_{\Delta}$  predicted by the different data-driven models for the two-dimensional shear-driven cavity flow at  $Re = 4250$ .

	Galerkin projection	Unc. quad. SINDy	Con. quad. SINDy	Unc. cubic SINDy	Con. cubic SINDy	N–S
$\sigma$	0.029	0.017	0.018	0.032	0.025	0.024
$\omega$	7.50	7.51	7.53	7.1	7.55	7.51
$\sigma_{\Delta}$	-0.11	-0.03	-0.035	-0.035	-0.03	N/A

TABLE 2. Comparison of growth rate  $\sigma$  and circular frequency  $\omega$  of the leading eigenvalue for the different models considered. Results from a global stability analysis of the Navier–Stokes equations are also reported. The last row also reports the value of the decay rate  $\sigma_{\Delta}$  associated with the shift mode.

on the determinant of the low-dimensional linear operator is however a non-convex constraint and does not fall within the scope of the library CVXOPT used in the present work.

Let us finally explore the decay rate  $\sigma_{\Delta}$  of the shift mode predicted by the different models, shown in table 2. Contrary to the cylinder flow, all models presented here correctly identify a linearly stable shift mode. Note however that if only the transient shown in figure 6(a) had been considered, the unconstrained quadratic model would suffer from the same shortcoming as for the cylinder flow (i.e. linearly unstable shift mode). Given that the relative distance to the critical Reynolds number is comparatively smaller in the present case compared to the cylinder flow, the influence of the truncated modes is expected to be less important. One might thus hypothesize that the correctness of the unconstrained quadratic model is related to this fact. In any case, this example once again underlines the importance of using physics-based constraints in order to identify physically relevant low-order models.

## 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. In fact, including higher-order

nonlinear terms in the models prevents underfitting, and allows for models with improved accuracy in terms of fewer, more energetic modes. The Galerkin regression framework leverages the recent sparse identification of nonlinear dynamics (SINDy) algorithm (Brunton *et al.* 2016a), and 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 two-dimensional cylinder flow and the shear-driven cavity have each been carefully analysed 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 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 some of the unconstrained models identified reproduce faithfully the dynamics of the original system, analysis of their spectral properties 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 can incorrectly identify 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 models which tend to under-estimate 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 optimization. Current developments, based on the nonlinear optimization 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.

Extending these constrained regression methods to experimental data may also present unique challenges and rewards. Many numerical schemes are designed to preserve energy; however, in turbulent simulations and experiments where sensors have limited bandwidth, dissipation may be large enough to affect the stability of models. Explicitly incorporating energy-preserving constraints in the SINDy regression may be especially important in these problems to find nearby conservative systems.

### Acknowledgements

We are grateful for many fruitful discussions with B. Noack, J. Proctor and N. Kutz. We also appreciate valuable feedback from S. Dawson and C. Rowley. S.L.B. acknowledges generous funding support from the US Defense Advanced Research Projects Agency (DARPA HR0011-16-C-0016) and from the US Air Force Office of Scientific Research (AFOSR FA9550-16-1-0650 and FA9550-18-1-0200).

### Appendix A. Connection with NARMAX

Over the years, a number of different approaches have been proposed for the identification of nonlinear dynamical systems from measured data. One of the most versatile and popular approach is NARMAX (Billings 2013): nonlinear auto-regressive model with exogeneous inputs. Although it targets the identification of discrete-time dynamical systems, its formulation is very close to that of the SINDy framework. Apart from the discrete-time versus continuous-time representations of the dynamics, one core difference between these two approaches essentially lies in the algorithm used to enforce the parsimony of the model: NARMAX classically uses the orthogonal least-squares procedure, while SINDy is based on a  $\ell_1$ -penalized or iterative hard-thresholded least-squares regression. In addition, the SINDy framework has been extended significantly, including to identify partial differential equations, to connect with Koopman operator theory, and to incorporate information criteria, rational function nonlinearities, and, in the present work, constraints. However, SINDy may be considered as a close relative of the NARMAX family of system identification, and the constrained SINDy algorithm may be used to identify effective NARMAX models, as demonstrated in this [Appendix A](#).

Here, we apply SINDy to identify NARMAX-like models of the two-dimensional cylinder flow at  $Re = 100$ ; we only consider this flow configuration for simplicity. As before, the first two transient evolutions depicted in figure 3(a,b) form the training dataset, while the last trajectory (see figure 3c) is used solely for cross-validation purposes. In the rest, we postulate that the discrete-time model can be written as

$$\mathbf{a}(t) = \mathcal{L}\mathbf{a}(t - \tau) + \mathcal{N}(\mathbf{a}(t - \tau)). \quad (\text{A } 1)$$

Given a library of admissible right-hand side functions  $\Theta(\mathbf{a})$ , this discrete-time system can be recast as

$$\mathbf{a}(t) = \Theta(\mathbf{a}(t - \tau))\mathbf{E}. \quad (\text{A } 2)$$

As before,  $\mathbf{E}_i$  is a sparse column vector indicating which functions from the library  $\Theta(\mathbf{a})$  are active in the equation governing the dynamics of  $a_i(t)$ . Given time-series data of  $\mathbf{a}(t)$ , these sparse column vectors can once again be identified using the SINDy algorithm, based on a  $\ell_1$ -regularized least-squares minimization. The last parameter



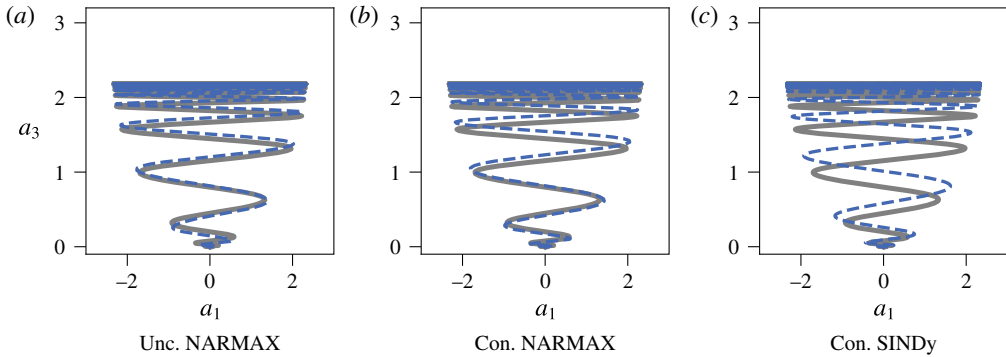


FIGURE 12. (Colour online) Comparison of the trajectory in the  $a_1$ – $a_3$  plane predicted by the different reduced-order models for the two-dimensional cylinder flow at  $Re = 100$ . The light grey trajectory is from direct numerical simulation.

which has to be set before the identification is the time lag  $\tau$ . Here, it is chosen as  $\tau = 0.25$ .

Figure 12 depicts the trajectory in the  $(a_1, a_3)$  plane predicted by an unconstrained NARMAX-like model, a constrained one and the constrained cubic Galerkin regression model presented earlier. Although all three trajectories are virtually identical, it must be noted that the two constrained models are more parsimonious than the unconstrained one. Indeed, while the unconstrained NARMAX model has 45 terms in its right-hand side, the constrained NARMAX and Galerkin regression models have 35 and 37 terms, respectively. Moreover, the unconstrained model fails to reproduce the transient evolution shown in figure 3(c). Importantly, both NARMAX models are identified using the SINDy algorithm and code base. This example clearly underlines the versatility of the SINDy framework and the benefit of introducing physics into the identification process by means of constraints.

## REFERENCES

- AKAIKE, H. 1974 A new look at the statistical model identification. *IEEE Trans. Autom. Control* **19** (6), 716–723.
- ÅKERVIK, E., BRANDT, L., HENNINGSON, D. S., HEPPFNER, J., MARXEN, O. & SCHLATTER, P. 2006 Steady solutions of the Navier–Stokes equations by selective frequency damping. *Phys. Fluids* **18** (6), 068102.
- ANDERSEN, M. S., DAHL, J. & VANDENBERGHE, L. 2013 CVXOPT: a Python package for convex optimization, version 1.1.6.
- BAGHERI, S. 2013 Koopman-mode decomposition of the cylinder wake. *J. Fluid Mech.* **726**, 596–623.
- BAGHERI, S., BRANDT, L. & HENNINGSON, D. S. 2009 Input–output analysis, model reduction and control of the flat-plate boundary layer. *J. Fluid Mech.* **620**, 263–298.
- BALAJEWICZ, M. J., DOWELL, E. H. & NOACK, B. R. 2013 Low-dimensional modelling of high-Reynolds-number shear flows incorporating constraints from the Navier–Stokes equation. *J. Fluid Mech.* **729**, 285–308.
- BARBAGALLO, A., SIPP, D. & SCHMID, P. J. 2009 Closed-loop control of an open cavity flow using reduced-order models. *J. Fluid Mech.* **641**, 1–50.
- BARKLEY, D. 2006 Linear analysis of the cylinder wake mean flow. *Europhys. Lett.* **75** (5), 750–756.
- BARKLEY, D. & HENDERSON, R. D. 1996 Three-dimensional Floquet stability analysis of the wake of a circular cylinder. *J. Fluid Mech.* **322**, 215–241.

- BERKOOZ, G., HOLMES, P. J. & LUMLEY, J. L. 1993 The proper orthogonal decomposition in the analysis of turbulent flows. *Annu. Rev. Fluid Mech.* **25** (1), 539–575.
- BILLINGS, S. A. 2013 *Nonlinear System Identification: NARMAX Methods in the Time, Frequency, and Spatio-Temporal Domains*. Wiley.
- BONGARD, J. & LIPSON, H. 2007 Automated reverse engineering of nonlinear dynamical systems. *Proc. Natl Acad. Sci. USA* **104** (24), 9943–9948.
- BRUNTON, S. L. & NOACK, B. R. 2015 Closed-loop turbulence control: progress and challenges. *Appl. Mech. Rev.* **67** (5), 050801.
- BRUNTON, S. L., PROCTOR, J. L. & KUTZ, J. N. 2016a Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proc. Natl Acad. Sci. USA* **113** (15), 3932–3937.
- BRUNTON, S. L., PROCTOR, J. L. & KUTZ, J. N. 2016b Sparse identification of nonlinear dynamics with control (SINDYc). *IFAC NOLCOS* **49** (18), 710–715.
- CANDÈS, E. J. 2006 Compressive sampling. In *Proceedings of the International Congress of Mathematicians*, vol. 3, pp. 1433–1452. European Mathematical Society.
- CARINI, M., AUTERI, F. & GIANNETTI, F. 2015 Centre-manifold reduction of bifurcating flows. *J. Fluid Mech.* **767**, 109–145.
- CARLBERG, K., BARONE, M. & ANTIL, H. 2017 Galerkin v. least-squares Petrov–Galerkin projection in nonlinear model reduction. *J. Comput. Phys.* **330**, 693–734.
- CARLBERG, K., TUMINARO, R. & BOGGS, P. 2015 Preserving Lagrangian structure in nonlinear model reduction with application to structural dynamics. *SIAM J. Sci. Comput.* **37** (2), B153–B184.
- CHARTRAND, R. 2011 Numerical differentiation of noisy, nonsmooth data. *ISRN Appl. Math.* **2011**, 164564.
- DONOHO, D. L. 2006 Compressed sensing. *IEEE Trans. Inform. Theory* **52** (4), 1289–1306.
- FABBIANE, N., SEMERARO, O., BAGHERI, S. & HENNINGSON, D. S. 2014 Adaptive and model-based control theory applied to convectively unstable flows. *Appl. Mech. Rev.* **66** (6), 060801.
- FISCHER, P. F., LOTTES, J. W. & KERKEMEIR, S. G. 2008 NEK5000: a fast and scalable high-order solver for computational fluid dynamics. <http://nek5000.mcs.anl.gov>.
- GLAZ, B., LIU, L. & FRIEDMANN, P. P. 2010 Reduced-order nonlinear unsteady aerodynamic modeling using a surrogate-based recurrence framework. *AIAA J.* **48** (10), 2418–2429.
- GLOERFELT, X. 2008 Compressible proper orthogonal decomposition/Galerkin reduced-order model of self-sustained oscillations in a cavity. *Phys. Fluids* **20** (11), 115105.
- GOLUB, G. H. & VAN LOAN, C. F. 2012 *Matrix Computations*, vol. 3. JHU Press.
- HAKEN, H. 1983 *Springer Series in Synergetics* (ed. M. Cardona, P. Fulde & H.-J. Queisser), p. 269. Springer.
- HOLMES, P. J., LUMLEY, J. L., BERKOOZ, G. & ROWLEY, C. W. 2012 *Turbulence, Coherent Structures, Dynamical Systems and Symmetry*, 2nd edn. Cambridge University Press.
- ILAK, M. & ROWLEY, C. W. 2008 Modeling of transitional channel flow using balanced proper orthogonal decomposition. *Phys. Fluids* **20**, 034103.
- ILLINGWORTH, S. J., MORGANS, A. S. & ROWLEY, C. W. 2010 Feedback control of flow resonances using balanced reduced-order models. *J. Sound Vib.* **330** (8), 1567–1581.
- JOHNSON, S. G. 2014 The NLOpt nonlinear-optimization package. <http://ab-initio.mit.edu/nlopt>.
- JONES, E., OLIPHANT, T., PETERSON, P. *et al.* 2001 SciPy: open source scientific tools for Python. <http://www.scipy.org/>.
- JUANG, J.-N. & PAPPAS, R. S. 1985 An eigensystem realization algorithm for modal parameter identification and model reduction. *J. Guid., Control Dyn.* **8** (5), 620–627.
- KAISER, E., NOACK, B. R., CORDIER, L., SPOHN, A., SEGOND, M., ABEL, M., DAVILLER, G., OSTH, J., KRAJNOVIC, S. & NIVEN, R. K. 2014 Cluster-based reduced-order modelling of a mixing layer. *J. Fluid Mech.* **754**, 365–414.
- KRIZHEVSKY, A., SUTSKEVER, I. & HINTON, G. E. 2012 Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems 25* (ed. F. Pereira, C. J. C. Burges, L. Bottou & K. Q. Weinberger), pp. 1097–1105. Curran Associates.
- KUKREJA, S. L. & BRENNER, M. J. 2007 Nonlinear system identification of aeroelastic systems: a structure-detection approach. In *Identification and Control*, pp. 117–145. Springer.

- KUKREJA, S. L., LÖFBERG, J. & BRENNER, M. J. 2006 A least absolute shrinkage and selection operator (lasso) for nonlinear system identification. *IFAC Proc.* **39** (1), 814–819.
- KUTZ, J. N. 2017 Deep learning in fluid dynamics. *J. Fluid Mech.* **814**, 1–4.
- KUTZ, J. N., BRUNTON, S. L., BRUNTON, B. W. & PROCTOR, J. L. 2016 *Dynamic Mode Decomposition: Data-Driven Modeling of Complex Systems*. SIAM.
- LEE, C., KIM, J., BABCOCK, D. & GOODMAN, R. 1997 Application of neural networks to turbulence control for drag reduction. *Phys. Fluids* **9** (6), 1740–1747.
- LING, J., KURZAWSKI, A. & TEMPLETON, J. 2016 Reynolds averaged turbulence modelling using deep neural networks with embedded invariance. *J. Fluid Mech.* **807**, 155–166.
- LINSCOTT, R. & WIKLUND, T. 2014 Parsimonious dynamical systems using the LASSO and the bootstrap. <http://uu.diva-portal.org/smash/get/diva2:750443/FULLTEXT01.pdf>.
- MAJDA, A. J. & HARLIM, J. 2012 Physics constrained nonlinear regression models for time series. *Nonlinearity* **26** (1), 201.
- MANGAN, N. M., BRUNTON, S. L., PROCTOR, J. L. & KUTZ, J. N. 2016 Inferring biological networks by sparse identification of nonlinear dynamics. *IEEE Trans. Mol. Biol. Multi-Scale Commun.* **2** (1), 52–63.
- MANGAN, N. M., KUTZ, J. N., BRUNTON, S. L. & PROCTOR, J. L. 2017 Model selection for dynamical systems via sparse regression and information criteria. *Proc. R. Soc. Lond. A* **473**, 20170009.
- MANTIČ-LUGO, V., ARRATIA, C. & GALLAIRE, F. 2014 Self-consistent mean flow description of the nonlinear saturation of the vortex shedding in the cylinder wake. *Phys. Rev. Lett.* **113** (8), 084501.
- MCCONAGHY, T. 2011 Ffx: Fast, scalable, deterministic symbolic regression technology. In *Genetic Programming Theory and Practice IX*, pp. 235–260. Springer.
- MELIGA, P. 2017 Harmonics generation and the mechanics of saturation in flow over an open cavity: a second-order self-consistent description. *J. Fluid Mech.* **826**, 503–521.
- MEZIĆ, I. 2005 Spectral properties of dynamical systems, model reduction and decompositions. *Nonlinear Dyn.* **41** (1–3), 309–325.
- MEZIĆ, I. 2013 Analysis of fluid flows via spectral properties of the Koopman operator. *Annu. Rev. Fluid Mech.* **45**, 357–378.
- MILANO, M. & KOUMOUTSAKOS, P. 2002 Neural network modeling for near wall turbulent flow. *J. Comput. Phys.* **182** (1), 1–26.
- NAIR, A. G. & TAIRA, K. 2015 Network-theoretic approach to sparsified discrete vortex dynamics. *J. Fluid Mech.* **768**, 549–571.
- NOACK, B. R., AFANASIEV, K., MORZYNSKI, M., TADMOR, G. & THIELE, F. 2003 A hierarchy of low-dimensional models for the transient and post-transient cylinder wake. *J. Fluid Mech.* **497**, 335–363.
- NOACK, B. R., MORZYNSKI, M. & TADMOR, G. 2011 *Reduced-Order Modelling for Flow Control*. vol. 528. Springer Science & Business Media.
- NOACK, B. R., STANKIEWICZ, W., MORZYNSKI, M. & SCHMID, P. J. 2016 Recursive dynamic mode decomposition of a transient cylinder wake. *J. Fluid Mech.* **809**, 843–872.
- ROSSITER, J. E. 1964 Wind tunnel experiments on the flow over rectangular cavities at subsonic and transonic speeds. *Tech. Rep.* Ministry of Aviation; Royal Aircraft Establishment; RAE Farnborough.
- ROWLEY, C. W. 2005 Model reduction for fluids using balanced proper orthogonal decomposition. *Intl J. Bifurcation Chaos* **15** (3), 997–1013.
- ROWLEY, C. W., COLONIUS, T. & BASU, A. J. 2002 On self-sustained oscillations in two-dimensional compressible flow over rectangular cavities. *J. Fluid Mech.* **455**, 315–346.
- ROWLEY, C. W. & DAWSON, S. 2017 Model reduction for flow analysis and control. *Annu. Rev. Fluid Mech.* **49**, 387–417.
- ROWLEY, C. W., MEZIĆ, I., BAGHERI, S., SCHLATTER, P. & HENNINGSON, D. S. 2009 Spectral analysis of nonlinear flows. *J. Fluid Mech.* **645**, 115–127.
- RUDY, S. H., BRUNTON, S. L., PROCTOR, J. L. & KUTZ, J. N. 2017 Data-driven discovery of partial differential equations. *Sci. Adv.* **3**, e1602614.
- SCHAEFFER, H. 2017 Learning partial differential equations via data discovery and sparse optimization. *Proc. R. Soc. Lond. A* **473**, 20160446.

- SCHLEGEL, M. & NOACK, B. R. 2015 On long-term boundedness of galerkin models. *J. Fluid Mech.* **765**, 325–352.
- SCHMID, P. J. 2010 Dynamic mode decomposition of numerical and experimental data. *J. Fluid Mech.* **656**, 5–28.
- SCHMIDT, M. & LIPSON, H. 2009 Distilling free-form natural laws from experimental data. *Science* **324** (5923), 81–85.
- SCHUMM, M., EBERHARD, B. & MONKEWITZ, P. A. 1994 Self-excited oscillations in the wake of two-dimensional bluff bodies and their control. *J. Fluid Mech.* **271**, 17–53.
- SCHWARZ, G. *et al.* 1978 Estimating the dimension of a model. *Ann. Stat.* **6** (2), 461–464.
- SEMAAN, R., KUMAR, P., BURNAZZI, M., TISSOT, G., CORDIER, L. & NOACK, B. R. 2016 Reduced-order modelling of the flow around a high-lift configuration with unsteady coanda blowing. *J. Fluid Mech.* **800**, 72–110.
- SEMERARO, O., LUSSEYRAN, F., PASTUR, L. & JORDAN, P. 2017 Qualitative dynamics of wavepackets in turbulent jets. *Phys. Rev. Fluids* **2**, 094605.
- SENGUPTA, T. K., HAIDER, S. I., PARVATHI, M. K. & PALLAVI, G. 2015 Enstrophy-based proper orthogonal decomposition for reduced-order modeling of flow past a cylinder. *Phys. Rev. E* **91** (4), 043303.
- SIPP, D. & LEBEDEV, A. 2007 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.
- SIPP, D., MARQUET, O., MELIGA, P. & BARBAGALLO, A. 2010 Dynamics and control of global instabilities in open-flows: a linearized approach. *Appl. Mech. Rev.* **63** (3), 030801.
- SIPP, D. & SCHMID, P. J. 2016 Linear closed-loop control of fluid instabilities and noise-induced perturbations: a review of approaches and tools. *Appl. Mech. Rev.* **68** (2), 020801.
- SIROVICH, L. 1987 Turbulence and the dynamics of coherent structures. Part I: coherent structures. *Q. Appl. Maths* **45** (3), 561–571.
- TADMOR, G., LEHMANN, O., NOACK, B. R. & MORZYŃSKI, M. 2010 Mean field representation of the natural and actuated cylinder wake. *Phys. Fluids* **22** (3), 034102.
- TIBSHIRANI, R. 1996 Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. B* **58**, 267–288.
- TU, J. H., ROWLEY, C. W., LUCHTENBURG, D. M., BRUNTON, S. L. & KUTZ, J. N. 2014 On dynamic mode decomposition: theory and applications. *J. Comput. Dyn.* **1** (2), 391–421.
- WANG, W. X., YANG, R., LAI, Y. C., KOVANIS, V. & GREBOGI, C. 2011 Predicting catastrophes in nonlinear dynamical systems by compressive sensing. *Phys. Rev. Lett.* **106**, 154101.
- WIGGINS, S. 2003 *Introduction to Applied Nonlinear Dynamical Systems and Chaos*, Texts in Applied Mathematics, vol. 2. Springer Science & Business Media.
- WILLCOX, K. & PERAIRE, J. 2002 Balanced model reduction via the proper orthogonal decomposition. *AIAA J.* **40** (11), 2323–2330.
- WILLIAMS, M. O., KEVREKIDIS, I. G. & ROWLEY, C. W. 2015 A data-driven approximation of the Koopman operator: extending dynamic mode decomposition. *J. Nonlinear Sci.* **25** (6), 1307–1346.
- YAMOUNI, S., SIPP, D. & JACQUIN, L. 2013 Interaction between feedback aeroacoustic and acoustic resonance mechanisms in a cavity flow: a global stability analysis. *J. Fluid Mech.* **717**, 134–165.
- YAO, C. & BOLLT, E. M. 2007 Modeling and nonlinear parameter estimation with Kronecker product representation for coupled oscillators and spatiotemporal systems. *Phys. D* **227** (1), 78–99.
- ZEBIB, A. 1987 Stability of viscous flow past a circular cylinder. *J. Engng Maths* **21** (2), 155–165.
- ZHANG, H.-Q., FEY, U., NOACK, B. R., KÖNIG, M. & ECKELMANN, H. 1995 On the transition of the cylinder wake. *Phys. Fluids* **7** (4), 779–794.
- ZHANG, W., WANG, B., YE, Z. & QUAN, J. 2012 Efficient method for limit cycle flutter analysis based on nonlinear aerodynamic reduced-order models. *AIAA J.* **50** (5), 1019–1028.
- ZHANG, Z. J. & DURAISAMY, K. 2015 Machine learning methods for data-driven turbulence modeling. In *22nd AIAA Computational Fluid Dynamics Conference*, p. 2460. American Institute of Aeronautics and Astronautics.