# ENC-2020-0646
# ADDRESSING OVERFITTING ISSUES IN THE
# SPARSE IDENTIFICATION OF NONLINEAR DYNAMICAL SYSTEMS

**Leo**nardo S. de B. **Alves**
Departamento de Engenharia Mecânica, Universidade Federal Fluminense, Niterói, RJ 24210-240, Brasil.
lsbalves@id.uff.br

*Abstract. Over the past four years, the derivation of dynamical models using symbolic regression has become a fixture in machine learning due to the development of the SINDy tool for the **S**parse **I**dentification of **N**onlinear **D**ynamical systems. As far as the author is aware, this tool has been applied assuming a polynomial representation of the unknown model that uses a monomial basis including only up to cubic nonlinearities. In the present paper, this issue is further explored. It turns out the library of candidate functions becomes ill-conditioned as the maximum nonlinearity order is increased, preventing the LASSO regularization in SINDy from eliminating unphysical terms due to increased error propagation.*

*Keywords: Machine Learning, Symbolic Regression, SINDy, Candidate Functions, Ill-Conditioned Library.*

## 1. INTRODUCTION

Data-based model development is at least as old as modern science, originally defined and championed by the Italian astronomer, physicist and engineer Galileu Galilei. Arguably its most well known pioneer is the German astronomer and mathematician Johannes Kepler, Galileu's contemporary, who used the planetary data he collected with his telescope to derive the now famous elliptical orbits to form his laws of planetary motion. This approach to scientific development, however, was not the one advocated by Galileu and, indeed, is not how the scientific method is generally understood today. This is in large part due to the impact of the work developed by English mathematician, physicist and astronomer Issac Newton, who was born around the time both passed away. He derived a dynamic model from first principles, describing the fundamental relationship between momentum and energy, that yielded the same elliptical orbits when applied to planetary motion. Experimental telescope data was employed only to validate the results obtained from this model.

Kepler's approach to model development is known today as symbolic regression. It is currently defined as a regression analysis that searches for the model that best fits an available dataset using a certain space of mathematical functions. Two recent milestones in the field of machine learning have put symbolic regression back at the forefront. The first one is related to the work by Bongard and Lipson (2007) as well as Schmidt and Lipson (2009), which identifies nonlinear differential equations from data using genetic programming (Koza, 1992). The second one is related to the work by Brunton *et al.* (2016), which takes advantage of sparse regression (Tibshirani, 1996) and compressed sensing (Donoho, 2006) to make symbolic regression dramatically less expensive. These tools have been applied to a wide range of applications, too large to be cited here. Fortunately, appropriate references do exist because this material has been reviewed quite recently. The reader is referred to the book by Brunton and Kutz (2018), which presents the fundamental theory required to understand and use these machine learning tools as well as a wide range of applications, and the review paper by Brunton *et al.* (2020), which does the same but in a summarized way and applied to problems in the field of fluid mechanics.

A far as the author is aware, most applications of this tool commonly known as SINDy (Brunton *et al.*, 2016), i.e. **S**parse **I**dentification of **N**onlinear **Dy**namics, consider only low order polynomial nonlinearities in the library of candidate functions as well as only problems that can be modeled by a low dimensional system of equations. For instance, this was the case when improving SINDy to include $i$) symmetry constraints in the development of reduced order models for the incompressible flow around a cylinder (Loiseau and Brunton, 2018), $ii$) the ability to discover hybrid dynamical systems, where the nonlinear governing equations change with time (Mangan *et al.*, 2019), and $iii$) guidelines for data sampling strategies and embeddings to improve the identification process (Champion *et al.*, 2019). The present paper explores these latter issues in further detail. It shows that the inclusion of high order nonlinearitie leads to the unavoidable increase in the condition number of the library of candidate functions, which prevents SINDy from identifying the correct nonlinear dynamical system due to extreme error propagation.

## 2. MATHEMATICAL FORMULATION

The discussion in the present paper will be restricted to systems of ordinary differential equations, since the long term goal of this work is the development of reduced order models derived using compressed data from projection techniques,

such as the Galerkin method, principal component analysis (PCA), proper orthogonal decomposition (POD) or dynamic mode decomposition (DMD). Hence, let us consider such a system in the form of

$$\frac{d\mathbf{x}}{dt} = \mathbf{f}\big(\mathbf{x}(t)\big) \quad , \tag{1}$$

where the state vector and function are written as

$$\mathbf{x}(t) = \big\{ x_1(t), x_2(t), \ldots, x_n(t) \big\}^T \quad \text{and} \quad \mathbf{f}\big(\mathbf{x}(t)\big) = \big\{ f_1\big(\mathbf{x}(t)\big), f_2\big(\mathbf{x}(t)\big), \ldots, f_n\big(\mathbf{x}(t)\big) \big\}^T \quad , \tag{2}$$

respectively.

## 2.1 The Tool

In order to use SINDy, a few assumptions must be made about Eqs. (1) and (2). They are:

1. The state size $n$ is arbitrary but small.

2. The state vector time history $\mathbf{x}(t)$ is available from data.

3. The state function $\mathbf{f}(\mathbf{x}(t))$ dependence on the state vector $\mathbf{x}(t)$ is unknown but sparse.

Although the first assumption does not affect the description of the tool, the second one is associated with the selection of the sampling rate $m$ and period $\tau = t_m - t_1$ required to build the matrixes

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}(t_1)^T \\ \mathbf{x}(t_2)^T \\ \vdots \\ \mathbf{x}(t_m)^T \end{pmatrix} = \begin{pmatrix} x_1(t_1) & x_2(t_1) & \cdots & x_n(t_1) \\ x_1(t_2) & x_2(t_2) & \cdots & x_n(t_2) \\ \vdots & \vdots & \ddots & \vdots \\ x_1(t_m) & x_2(t_m) & \cdots & x_n(t_m) \end{pmatrix}_{m \times n} \quad \text{and} \tag{3}$$

$$\dot{\mathbf{X}} = \begin{pmatrix} \dot{\mathbf{x}}(t_1)^T \\ \dot{\mathbf{x}}(t_2)^T \\ \vdots \\ \dot{\mathbf{x}}(t_m)^T \end{pmatrix} = \begin{pmatrix} \dot{x}_1(t_1) & \dot{x}_2(t_1) & \cdots & \dot{x}_n(t_1) \\ \dot{x}_1(t_2) & \dot{x}_2(t_2) & \cdots & \dot{x}_n(t_2) \\ \vdots & \vdots & \ddots & \vdots \\ \dot{x}_1(t_m) & \dot{x}_2(t_m) & \cdots & \dot{x}_n(t_m) \end{pmatrix}_{m \times n} \quad , \tag{4}$$

where the latter matrix can be either measured directly or approximated numerically from the former matrix. Due to the third assumption, one must now construct a library of candidate functions to estimate the unknown dependence of $\mathbf{f}(\mathbf{x}(t))$ on $\mathbf{x}(t)$. If no information is available from the problem under investigation, a polynomial representation using a monomial basis is often the preferred choice. The library of candidate functions $\mathbf{F}$ can then be defined as

$$\mathbf{F} = \begin{pmatrix} 1 & x_1(t_1) & x_2(t_1) & \cdots & x_n(t_1) & x_1(t_1)^2 & x_1(t_1)\,x_2(t_1) & \cdots & x_n(t_1)^2 & \cdots & x_n(t_1)^p \\ 1 & x_1(t_2) & x_2(t_2) & \cdots & x_n(t_2) & x_1(t_2)^2 & x_1(t_2)\,x_2(t_2) & \cdots & x_n(t_2)^2 & \cdots & x_n(t_2)^p \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_1(t_m) & x_2(t_m) & \cdots & x_n(t_m) & x_1(t_m)^2 & x_1(t_m)\,x_2(t_m) & \cdots & x_n(t_m)^2 & \cdots & x_n(t_1m)^p \end{pmatrix}_{m \times q} \tag{5}$$

where $p$ is the maximum chosen nonlinearity order for the resulting polynomial representation and $q$ is the total number of terms in the library. Hence, defining the unknown linear coefficient matrix $\mathbf{C}$ as

$$\mathbf{C} = \begin{pmatrix} c_{1,1} & c_{1,2} & \cdots & c_{1,n} \\ c_{2,1} & c_{2,2} & \cdots & c_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ c_{q,1} & c_{q,2} & \cdots & c_{q,n} \end{pmatrix}_{q \times n} \quad , \tag{6}$$

means one can estimate the functional form of each $f_j(\mathbf{x}(t))$ from its respective data set $\mathbf{X}_j$ by solving

$$\dot{\mathbf{X}}_j^T = \mathbf{F} \cdot \mathbf{C}_j^T \quad , \tag{7}$$

instead of Eq. (1) for each column $j = 1, 2, \ldots, n$.

It is important to note that $m \gg q$ in most applications, i.e. $\mathbf{F}$ is a rectangular (low rank) matrix. Hence, Eq. (7) is over-determined and, in fact, represents a linear optimization problem. In order to take advantage of the third assumption, the objective function $\mathbf{u}_j$ to be minimized is usually defined as

$$\mathbf{u}_j = ||\,\dot{\mathbf{X}}_j^T - \mathbf{F} \cdot \mathbf{C}_j^T\,||_2 + \lambda\,||\,\mathbf{C}_j^T\,||_1 \tag{8}$$

where $\lambda$ is the sparsity identification LASSO regularization parameter (Tibshirani, 1996).

## 2.2 The Test Case

SINDy's performance is evaluated in the present paper by considering the Lorenz equations (Sparrow, 1982) as the model problem that provides the data set. These equations take the form

$$
\begin{aligned}
\dot{x}(t) &= \sigma\left(y(t) - x(t)\right) \quad, \\
\dot{y}(t) &= x(t)\left(\rho - z(t)\right) - y(t) \quad \text{and} \\
\dot{z}(t) &= x(t)\,y(t) - \beta\,z(t) \quad,
\end{aligned} \tag{9}
$$

where $\sigma$, $\rho$ and $\beta$ are the model parameters. Three different parametric conditions are considered here. Although all of them provide time series data during the time asymptotic nonlinear behavior range, each behavior is quite different from one another. Their disturbance (left column) temporal behavior and (right column) spectra are shown in Fig. 1, representing (top row) chaotic, (middle row) double periodic and (bottom row) periodic asymptotic nonlinear regimes. In this figure, $x_S$ represents one of the two nontrivial steady-states.
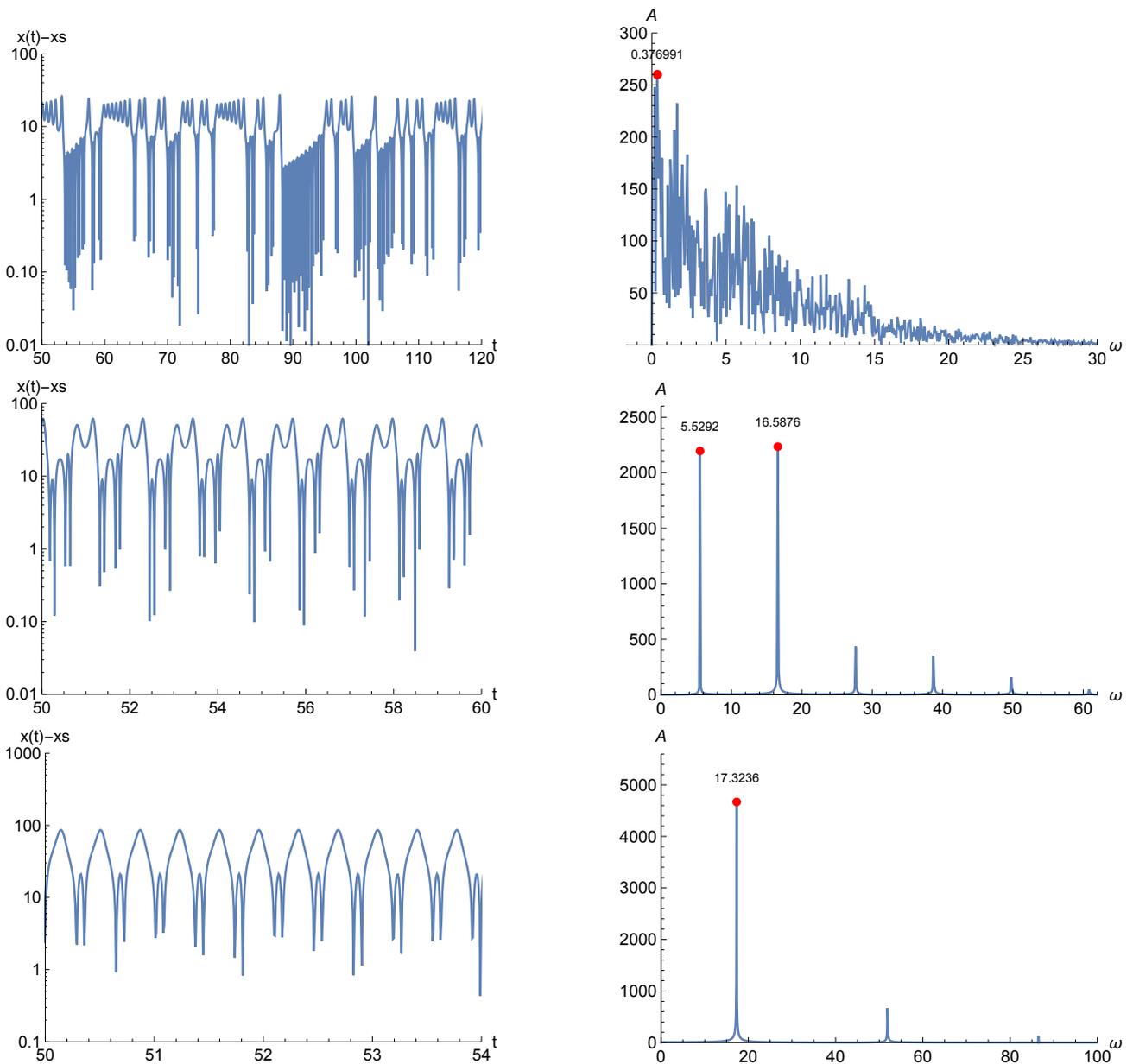


Figure 1. Disturbance (left column) temporal behavior and (right column) spectra associated with the $x(t)$ solution from the Lorenz equations provided in Eq. (9) obtained with $\sigma = 10$ and $\beta = 8/3$, which yields (top row) a chaotic regime when $\rho = 28$, (middle row) a double periodic regime when $\rho = 165$ and (bottom row) a periodic regime when $\rho = 400$.

## 3. RESULTS

The first case considered here assumes an inverse problem (Özisik and Orlande, 2000) scenario instead of a typical machine learning one, where the correct model is known *a priori*. In this particular case, all seven linear coefficients are considered unknown instead of just the ones associated with the three unknown parameters. Very accurate values of the latter were obtained even without regularization, i.e. $\lambda = 0$. For instance, they are $\sigma = 10.000000$, $\beta = 2.6666667$ and $\rho = 165.00000$ under the double periodic parametric conditions. However, these errors increase when the maximum nonlinearity order chosen also increases, e.g. $(\sigma, \beta, \rho) = (10.000000, 2.6666682, 165.00000)$, $(10.000001, 2.6666702, 165.00000)$, $(10.000005, 2.6668199, 165.00001)$ and $(10.034239, 3.4425656, 165.17240)$ as $p$ is increased from 2 to 5, respectively. Furthermore, it becomes increasingly more difficult to distinguish which terms in the estimated model are not part of the original model. This is illustrated in Fig. 2, which shows all coefficients for each equation normalized by its respective maximum coefficient, for $p = 5$, $t_1 = 30$, $t_m = 37.575758$ and $m = 500$.
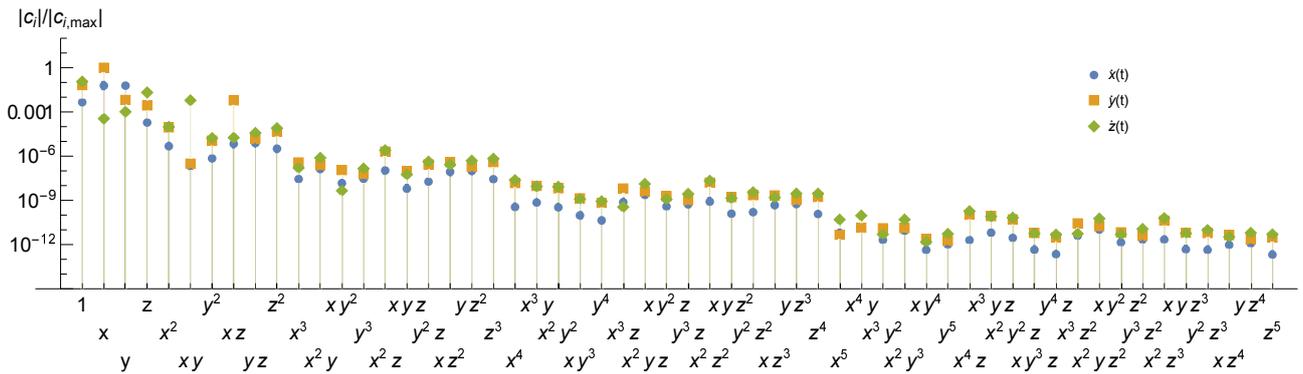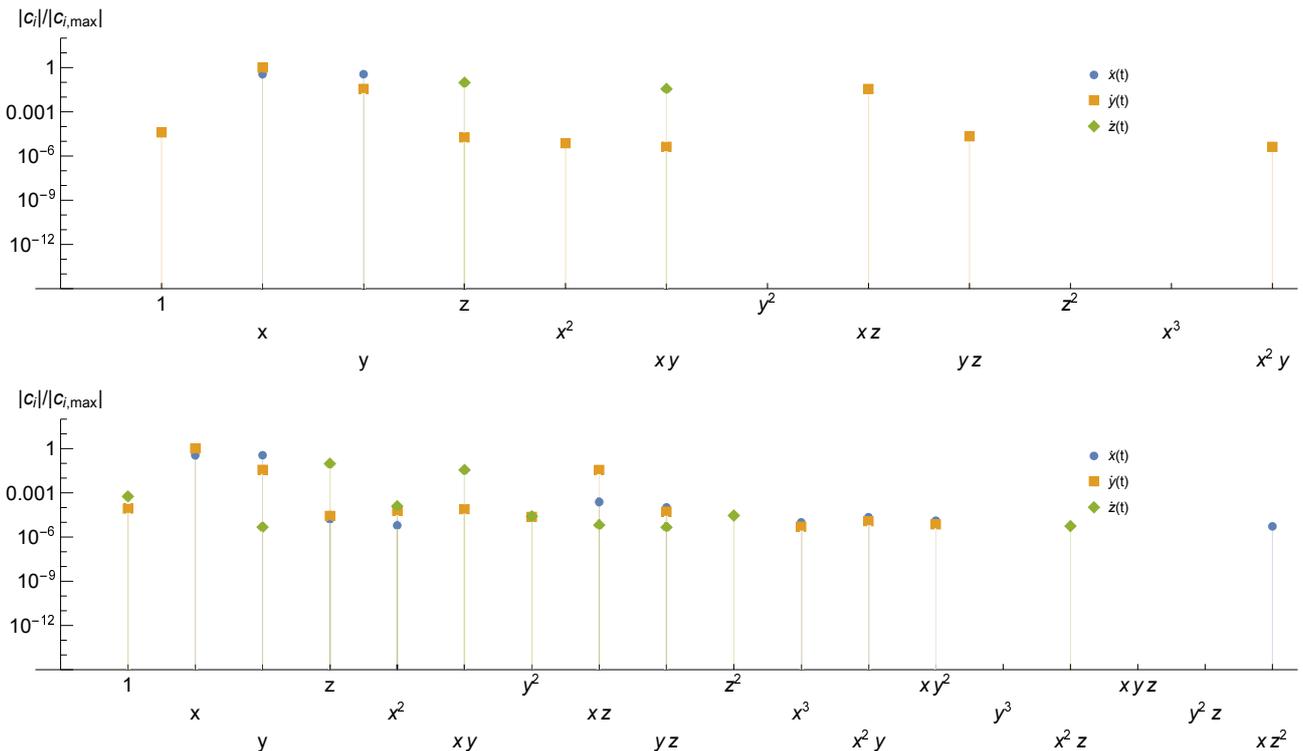


Figure 2. Normalized linear coefficients for the double periodic parametric conditions obtained by SINDy without a LASSO regularization ($\lambda = 0$) with $p = 5$, $t_1 = 30$, $t_m = 37.575758$ and $m = 500$.

Unfortunately, adding a LASSO regularization does not fix this problem. Figure 3 presents the normalized linear coefficients obtained under chaotic parametric conditions for $p = 4$, $t_1 = 50$, $t_m = 150$ and $m = 500$ when $\lambda = 10^{-1}$ (top), $10^{-2}$ (middle) and $10^{-4}$ (bottom). A significant number of unphysical terms could not be removed from the estimated model independent of the regularization parameter magnitude. Although not shown here, using $\lambda > 10^{-1}$ leads to similar results and $\lambda < 10^{-4}$ leads to the same plot shown at the bottom of this figure.
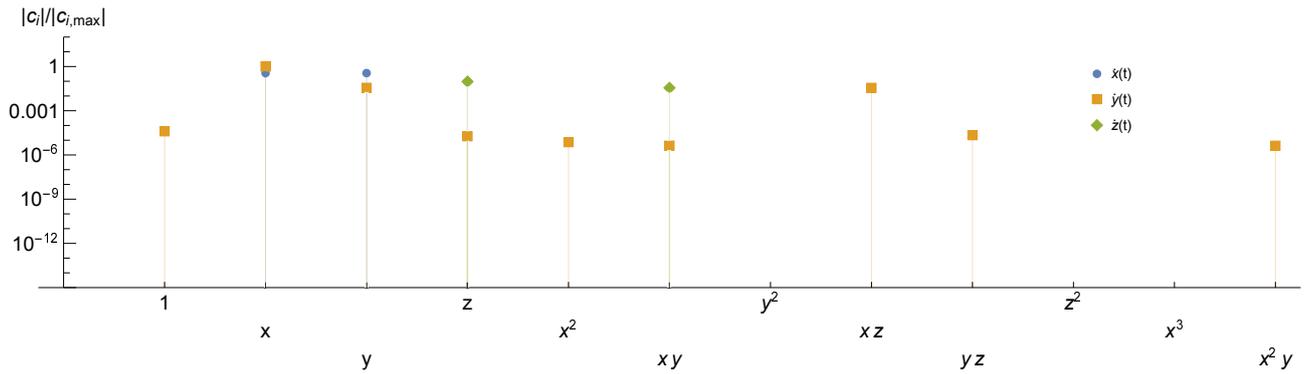
Figure 3. Normalized linear coefficients for the chaotic parametric conditions obtained by SINDy for $p = 4$, $t_1 = 50$, $t_m = 150$ and $m = 500$ using a LASSO regularization with $\lambda = 10^{-1}$ (top), $10^{-2}$ (middle) and $10^{-4}$ (bottom).

A numerical analysis of the error bounds associated with the solution of a linear system in the form of $\mathbf{b} = \mathbf{A} \cdot \mathbf{x}$ reveals that the condition number $\kappa_i(\mathbf{A})$ of matrix $\mathbf{A}$ is a measure of the sensitivity of the unknown vector $\mathbf{x}$ to errors in both the matrix $\mathbf{A}$ and the known vector $\mathbf{b}$ (Stoer and Bulirsch, 1993). The condition number is defined as $\kappa_i(\mathbf{A}) = ||\mathbf{A}||_i \, ||\mathbf{A}^{-1}||_i$, where $|| \cdot ||_i$ represents an arbitrary norm $i$ and $\mathbf{A}^{-1}$ is the inverse of $\mathbf{A}$. It is important to notice, however, that this statement is rigorously true for square matrixes. Nevertheless, it is also valid for rectangular linear systems in the form of Eq. (7) as long as two additional restrictions are imposed: $i$) a pseudo-inverse is employed instead of a regular inverse and $ii$) either a maximum or an Euclidian norm is employed instead of an arbitrary norm (Demko, 1986). In the present work, we use $\kappa_2(\mathbf{F}) = ||\mathbf{F}||_2 \, ||\mathbf{F}^{\dagger}||_2$, where $\mathbf{F}^{\dagger}$ is the Moore-Penrose inverse of $\mathbf{F}$ and $|| \cdot ||_2$ represents an Euclidian norm. This relationship between condition number (left) and relative error (right) is shown in Fig. 4 for the double periodic parametric conditions, both as functions of the sampling rate $m$ for different sampling periods $\tau$, where $T$ is the period of the dominant mode. First and foremost, this figure confirms that the condition number is an adequate proxy for the relative error, as long as a minimum sampling rate is employed. If smaller values are used, there is not enough data available for the linear regression to work appropriately. Furthermore, there is a maximum sampling rate and a maximum sampling period beyond which neither condition number nor relative error decreases. These results indicate that the development of an accurate ROM becomes significantly more efficient when the condition number $\kappa_2(\mathbf{F})$ can be calculated and used to identify the optimal sampling rate and period *a priori*, improving known sampling strategies (Champion *et al.*, 2019).
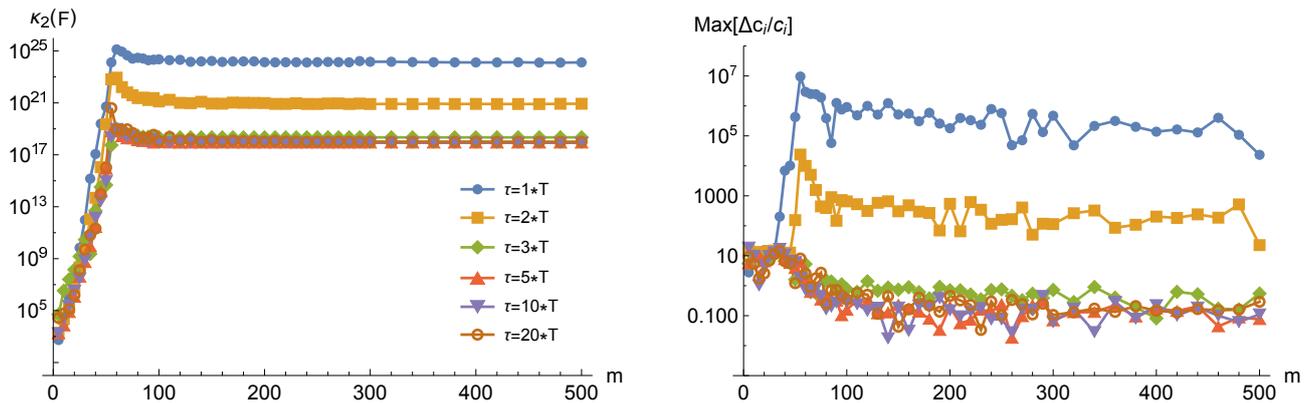


Figure 4. Condition number of the library matrix (left) and maximum linear coefficient relative error (right) for the double periodic parametric conditions with $p = 5$, $t_1 = 30$ and $T = 0.378788$ without the LASSO regularization.

Although not shown here, the information presented in Fig. 4 was generated for the other two parametric conditions as well. It was then used to extract optimal sampling rates and periods for all three parametric conditions with different maximum nonlinearity orders. Optimal here means minimizing condition number and/or relative error. Furthermore, nonlinearity orders were varied from $p = 1$ to 5, where 1 represents the original model. These results for the optimal (left) condition number and (right) maximum relative error are presented in Fig. 5 as functions of $p$. The first important result from this figure is that both increase as the maximum nonlinearity order increases. This is not an entirely unexpected result if one realizes that the library of candidate functions $\mathbf{F}$ assumes the form of a Vandermonde matrix when the polynomial representation of these functions employ a monomial basis. Vandermonde matrixes are badly ill-conditioned matrixes that become more ill-conditioned as their size increases (Pan, 2016). This explains why the LASSO regularization fails
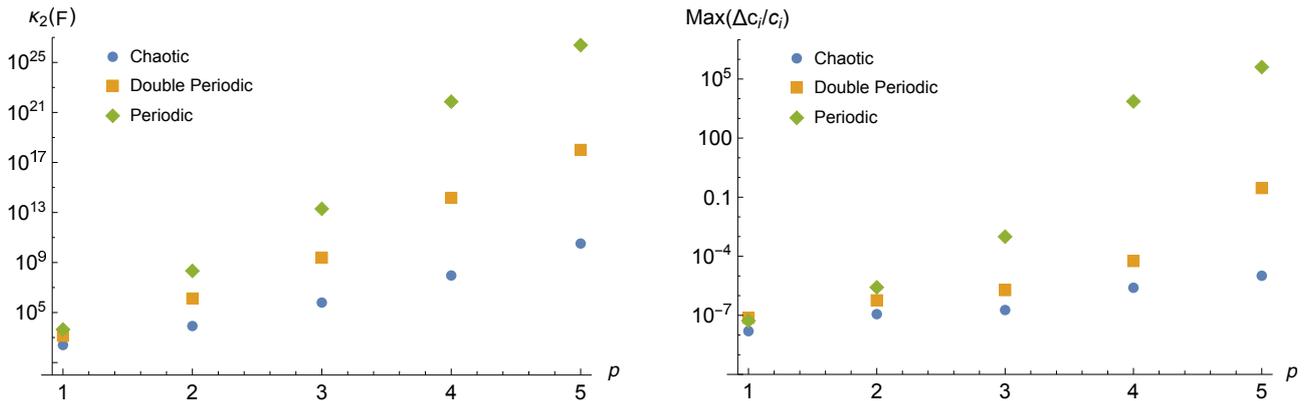
Figure 5. Optimal (left) condition number of the library matrix and (right) maximum linear coefficient relative error for all three parametric conditions as functions of maximum nonlinearity order without the LASSO regularization.

as the maximum nonlinearity order increases. It introduces a penalty for small coefficients in the objective function. Doing so, however, also introduces more error. Its magnitude is then amplified when solving the linear system in Eq. (7), increasingly so as the library condition number increases. This result also indicates that increasing the state size $n$ leads to similar error propagation issues, since doing so increases the size of the library matrix.

There are a few exceptions though, such as the Vandermonde matrixes arising from a discrete Fourier transform. The key property of a low condition number matrix is having its knots equally spaced around the unit circle. In the DFT case, they are all unitary up to scaling, which guarantees a perfect condition number, i.e. $\kappa = 1$, for square matrixes. This might explain the second important result in Fig. 5, which is the somewhat counterintuitive fact that the both condition number and maximum relative error are smallest for the chaotic parametric conditions, followed by the double periodic and periodic parametric conditions, when $2 \leq p \leq 5$. Periodicity in whatever form likely makes the knots unequally spaced around the unit circle. Small deviations from the equally spaced condition leads to large increases in condition number (Pan, 2016). On the other hand, chaotic conditions likely induce a random coefficient magnitude distribution inside the library matrix. Matrixes with such structures are known to have lower condition numbers (Edelman, 1988).

## 4. CONCLUSIONS AND FUTURE WORK

The above discussion clearly indicates that the Vandermonde nature of the library matrix of candidate functions leads to an increasingly more ill-conditioned linear regression problem for the unknown coefficient as the maximum nonlinearity order and/or state size increase, which severely restricts the usability of SINDy. This is caused by the monomial basis employed in the polynomial representation of the unknown nonlinear dynamical system.

Hence, a path forward in an attempt to solve or minimize this issue is to change the basis employed. Different orthogonal bases are currently being tested and these results will be presented in the future.

## 5. ACKNOWLEDGEMENTS

## 6. REFERENCES

Bongard, J. and Lipson, H., 2007. "Automated reverse engineering of nonlinear dynamical systems". *Proceedings of the National Academy of Sciences*, Vol. 104, No. 24, pp. 9943–9948.

Brunton, S.L. and Kutz, J.N., 2018. *Data-Driven Science and Engineering: Machine Learning, Dynamical Systems, and Control*. Cambridge University Press.

Brunton, S.L., Noack, B.R. and Koumoutsakos, P., 2020. "Machine learning for fluid mechanics". *Annual Review of Fluid Mechanics*, Vol. 52, pp. 477–508.

Brunton, S.L., Proctor, J.L. and Kutz, J.N., 2016. "Discovering governing equations from data by sparse identification of nonlinear dynamical systems". *Proceedings of the National Academy of Sciences*, Vol. 113, No. 15, pp. 3922–3937.

Champion, K.P., Brunton, S.L. and Kutz, J.N., 2019. "Discovery of nonlinear multiscale systems: Sampling strategies and embeddings". *SIAM Journal of Applied Dynamical Systems*, Vol. 18, No. 1, pp. 312–333.

Demko, S., 1986. "Condition numbers of rectangular systems and bounds for generalized inverses". *Linear Algebra and Its Applications*, Vol. 78, pp. 199–206.

Donoho, D.L., 2006. "Compressed sensing". *IEEE Transactions on Information Theory*, Vol. 52, No. 4, pp. 1289–1306.

Edelman, A., 1988. "Eigenvalues and condition numbers of random matrices". *SIAM Journal on Matrix Analysis and Applications*, Vol. 9, No. 4, p. 8.

Koza, J.R., 1992. *Genetic Programming: On the Programming of Computers by Means of Natural Selection*, Vol. 1. MIT Press.

Loiseau, J.C. and Brunton, S.L., 2018. "Constrained sparse Galerkin regression". *Journal of Fluid Mechanics*, Vol. 838, pp. 42–67.

Mangan, N.M., Askham, T., S. L. Brunton and, J.N.K. and Proctor, J.L., 2019. "Model selection for hybrid dynamical systems via sparse regression". *Proceedings of the Royal Society of London*, Vol. 475, No. 20180534.

Özisik, M.N. and Orlande, H.R.B., 2000. *Inverse Heat Transfer: Fundamentals and Applications*. Taylot & Francis.

Pan, V.Y., 2016. "How bad are Vandermonde matrices?" *SIAM Journal on Matrix Analysis and Applications*, Vol. 37, No. 2, pp. 676–694.

Schmidt, M. and Lipson, H., 2009. "Distilling free-form natural laws from experimental data". *Science*, Vol. 324, pp. 81–85.

Sparrow, C., 1982. *The Lorenz Equations: Bifurcations, Chaos and Strange Attractors*, Vol. 41 of *Applied Mathematical Sciences*. Springer-Verlag, New York.

Stoer, J. and Bulirsch, R., 1993. *Introduction to Numerical Analysis*. Springer – Verlag, New York, 2nd edition.

Tibshirani, R., 1996. "Regression shrinkage and selection via the lasso". *Journal of the Royal Statistical Society. Series B.*, Vol. 58, No. 1, pp. 267–288.

## 7. RESPONSIBILITY NOTICE

The following text, properly adapted to the number of authors, must be included in the last section of the paper:
The author(s) is (are) solely responsible for the printed material included in this paper.