

Data-driven Discovery of Closure Models*

Shaowu Pan[†] and Karthik Duraisamy[†]

Abstract. Derivation of reduced order representations of dynamical systems requires the modeling of the truncated dynamics on the retained dynamics. In its most general form, this so-called closure model has to account for memory effects. In this work, we present a framework of operator inference to extract the governing dynamics of closure from data in a compact, non-Markovian form. We employ sparse polynomial regression and artificial neural networks to extract the underlying operator. For a special class of non-linear systems, observability of the closure in terms of the resolved dynamics is analyzed and theoretical results are presented on the compactness of the memory. The proposed framework is evaluated on examples consisting of linear to nonlinear systems with and without chaotic dynamics, with an emphasis on predictive performance on unseen data.

Key words. data-driven closures, dynamical system closures, reduced order modeling, machine learning

AMS subject classifications. 70G60, 76F20

1. Introduction. Complex problems in science and engineering are typically characterized by high-dimensional dynamics. Examples include the modeling of turbulent fluid flows, molecular dynamics, and astrophysical plasmas. When such problems are viewed from a dynamical systems perspective, the high dimensionality of phase space is a consequence of the fact that important physical processes occur over a wide range of spatial and temporal scales. However, effective computational models of these systems for the purposes of analysis, design and control require accurate low-dimensional representations. Popular techniques to obtain low-dimensional representations include projection-based reduced order models [22, 6, 52, 9], reduced basis methods [50, 37], proper generalized decomposition [10], and Krylov subspace techniques [4]. All of these techniques aim to capture the dynamics essential to a quantity of interest in by solving for a small number of unknowns (usually by restricting the dynamics to a low-dimensional manifold). In most practical situations, however, the multiscale nature of the problem is such that a low-dimensional representation requires closure. In other words, the influence of the discarded degrees of freedom on the retained unknowns becomes important and must be modeled.

The closure problem is well-recognized by the scientific computing community, and is typically addressed by invoking physical and/or mathematical arguments. A pertinent example of physics-based closure is Large Eddy Simulation [38] (LES) in fluid dynamics, where the impact of the unresolved scales on the resolved scales is often modeled via an eddy diffusivity hypothesis [18]. Another example [27] involves the determination of constitutive properties of complex materials through the detailed modeling of the microstructure. Approximate Green's function-based closures [24], adaptive deconvolution [43], and homogenization techniques [33] are representative of mathematically-inspired closures.

*Submitted to the editors 03/25/2018.

Funding: AFOSR grants FA9550-16-1-0309 & FA9550017-1-0195

[†]Department of Aerospace Engineering, University of Michigan, Ann Arbor, MI (shawnpa@umich.edu, kdur@umich.edu)

An alternate approach is to pursue data-driven techniques to address closure. There are several instances of the use of data in closure modeling and the following is not intended to be a complete or chronological review, but rather presents representative references from the viewpoint of the various levels at which data has been used to aid closure modeling. Observational data has been used to calibrate closure parameters in reduced fidelity [31] and reduced order [14, 5] models. In these approaches, the functional form of the closure term is prescribed (for instance, via an eddy viscosity assumption) with free parameters which are inferred by minimizing the misfit between the model output and training data. As an example of a more extensive approach, Xie et al. [54] impose a general structure to the closure term and infer matrix operators within the structure. At the next higher level, Ibanez et al. [25] use manifold learning to identify locally-linear embeddings and construct constitutive relationships for elasticity. Parish et al. [32], Singh et al. [42] directly extract the functional form of augmentations to the closure term by combining statistical inference and learning.

The goal of this work is to extract closure operators for reduced-dimensional dynamical systems using data snapshots generated from the original high-dimensional dynamical system. The low-dimensional state is augmented with a new set of variables, which represent the closure term, and the evolution equation for the dynamics is discovered in terms of the low-dimensional state using polynomial regression and neural networks. A key difference compared to the literature cited above is that the closed lower-dimensional system is capable of emulating non-Markovian characteristics. Further, the functional form of the evolution equation of the closure is not imposed, but rather extracted directly from the data.

Over the past few years, much research has been dedicated to the topic of “data-driven discovery of governing equations,” using techniques such as dynamic mode decomposition [40], feature-space regression [53, 8], operator inference [35], and neural networks [36], etc. These works have demonstrated that it is possible to a) *rediscover* known equations from data, or b) derive approximate representations of systems for which precise equations cannot be written (such as the spread of epidemics [29]). The scope of the present work is different, as the structure of the closure is unknown even for simple non-linear dynamical systems. It is, however, assumed that the full-order model corresponding to the high-dimensional system is known (as is the case in many physical problems, such as fluid dynamics where the governing equations are known, but are prohibitively expensive to solve in high-dimensional form) and this knowledge is incorporated into the model formulation process. Furthermore, emphasis is on prediction rather than reconstruction.

This paper is structured as follows: In [section 2](#), the closure problem is briefly described. In [section 3](#), a framework of operator inference is presented. In [section 4](#) and [section 6](#), applications of this framework with sparse polynomial regression and artificial neural network (ANN) are presented on various problems ranging from simple linear systems to a nonlinear PDE system. Theoretical investigations are conducted on the structure of the closure dynamics in [section 5](#). Conclusions and perspectives are drawn in [section 7](#).

2. Description of the closure problem. Consider the autonomous nonlinear dynamical system in (2.1)

$$(2.1) \quad \frac{d\mathbf{x}}{dt} = \mathbf{F}(\mathbf{x}),$$

where $\mathbf{x}(t) \in \mathbb{R}^N$, $N \in \mathbb{N}^+$, $t \in [0, +\infty)$, $\mathbf{x}(0) = \mathbf{x}_0$ and $\mathbf{F}(\cdot) : \mathbb{R}^N \mapsto \mathbb{R}^N$.

To serve as a representative lower-dimensional dynamical system, we consider a partition

$$(2.2) \quad \mathbf{x} = \begin{bmatrix} \hat{\mathbf{x}} \\ \tilde{\mathbf{x}} \end{bmatrix}, \mathbf{F}(\mathbf{x}) = \begin{bmatrix} \hat{\mathbf{F}}(\hat{\mathbf{x}}, \tilde{\mathbf{x}}) \\ \tilde{\mathbf{F}}(\hat{\mathbf{x}}, \tilde{\mathbf{x}}) \end{bmatrix},$$

where $\hat{\mathbf{F}}(\cdot, \cdot) : \mathbb{R}^Q \times \mathbb{R}^{N-Q} \mapsto \mathbb{R}^Q$, $\tilde{\mathbf{F}}(\cdot, \cdot) : \mathbb{R}^Q \times \mathbb{R}^{N-Q} \mapsto \mathbb{R}^{N-Q}$, $Q \in \mathbb{N}^+$. $\hat{\mathbf{x}} \in \mathbb{R}^Q$ is the low-dimensional or resolved state and $\tilde{\mathbf{x}} \in \mathbb{R}^{N-Q}$ represents the unresolved modes. In general terms, the above partition appears arbitrary. This partition is, however, directly relevant in a number of problems: (i) in projection-based Reduced Order Models (ROMs), where the components of the state in the original dynamical system are ordered according to an energy metric; (ii) in large eddy simulations (LES) of turbulence using spectral or finite element methods, where there is a clear separation between resolved and unresolved scales; (iii) in system identification, where the system is only partially observed and a governing equation for the partially observed system is desired.

The evolution of the reduced state is given by

$$(2.3) \quad \frac{d\hat{\mathbf{x}}}{dt} = \hat{\mathbf{F}}(\hat{\mathbf{x}}, \tilde{\mathbf{x}}),$$

where $\hat{\mathbf{x}}(0) = \hat{\mathbf{x}}_0 \in \mathbb{R}^Q$ and $\tilde{\mathbf{x}}(0) = \tilde{\mathbf{x}}_0 \in \mathbb{R}^{N-Q}$.

Note that (2.3) is not very useful, as the trajectory of the unresolved state $\tilde{\mathbf{x}}(t)$ is present in these equations. In reduced order modeling, a closed set of equations of the form

$$(2.4) \quad \frac{d\hat{\mathbf{x}}}{dt} = \mathbf{F}_{ROM}(\hat{\mathbf{x}}),$$

is obtained through physically-inspired [14, 5], data-augmented [55] or purely data-driven [35] methods.

In 2.4, $\hat{\mathbf{x}}(0) = \hat{\mathbf{x}}_0$ and $\mathbf{F}_{ROM} : \mathbb{R}^Q \mapsto \mathbb{R}^Q$. When $\mathbf{F}_{ROM}(\hat{\mathbf{x}}) = \hat{\mathbf{F}}(\hat{\mathbf{x}}, \mathbf{0})$, one obtains a classic truncated ROM. In fluid dynamic modeling terms, this corresponds to a Large Eddy Simulation without a explicit subgrid scale model. In obtaining such approximations, a key fact to consider is that, even if the high-dimensional system is Markovian, the corresponding projected low-dimensional system can be non-Markovian. This is evident even in the simplest case of a linear system. Consider that the full order model with its partition into resolved and unresolved states:

$$(2.5) \quad \frac{d}{dt} \begin{bmatrix} \hat{\mathbf{x}} \\ \tilde{\mathbf{x}} \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{x}} \\ \tilde{\mathbf{x}} \end{bmatrix}.$$

The evolution of the resolved states is given by:

$$(2.6) \quad \frac{d\hat{\mathbf{x}}}{dt} = A_{11}\hat{\mathbf{x}} + \int_0^t A_{12}e^{A_{22}\tau}A_{21}\hat{\mathbf{x}}(t-\tau)d\tau + A_{12}e^{A_{22}t}\tilde{\mathbf{x}}(0).$$

This equation is closed in the resolved variables. The first term on the right hand side is $\hat{\mathbf{F}}(\hat{\mathbf{x}}, \mathbf{0})$; the second term, which represents the closure, involves the time-history of the

resolved modes, and the third term involves the initial condition of the unresolved state, and can be expected to decay in time for a dissipative system.

The above expression of closure can be generalized to any nonlinear system (2.1) and (2.2) using the Mori-Zwanzig formalism [13, 12], and exact evolution equations can be written for the reduced state in the following Generalized Langevin form:

$$(2.7) \quad \frac{d\hat{\mathbf{x}}}{dt} = \hat{\mathbf{F}}(\hat{\mathbf{x}}, \mathbf{0}) + \int_0^t \mathcal{K}(\hat{\mathbf{x}}(s), t-s) ds + \mathcal{Q}(\hat{\mathbf{x}}(t)),$$

where \mathcal{K} , \mathcal{Q} are complex operators associated with convolution and influence of the initial conditions, respectively.

While the above equation is mathematically precise and formally closed in the resolved variables, the functional forms of \mathcal{K} and \mathcal{Q} are not explicitly known and numerically intractable, even for the simplest non-linear dynamical systems, and must thus be determined via the solution of another high-dimensional partial differential equation [20]. The key message is that reduced-order representations of even a linear Markovian system can introduce memory or time-history effects in an explicit form that requires *all* its previous states. In the present work, we introduce a dynamic memory and aim to extract its evolution using operator inference. It is also shown that for a specialized class of nonlinear systems that the memory length is compact, and thus the full history of resolved states is not necessary.

3. Framework of operator inference. As indicated by the Mori-Zwanzig formalism, the exact closure is a compositional convolution operator on all past resolved states. This approach is equivalent to the concept of dynamic or recurrent memory [19], a concept which has been very attractive in the time series modeling and deep learning communities. To address complex memory structures, we consider time delay vectors in the framework as implied by Takens embedding theorem [45] which states that there exists a diffeomorphism between proper time delayed (reconstructed) attractor and the original manifold.

By leveraging both dynamic memory and implications of Takens embedding theorem, a framework of operator inference is proposed as shown in the following discretized augmented system:

$$(3.1) \quad \frac{D\hat{\mathbf{x}}}{Dt} = \hat{\mathbf{F}}(\hat{\mathbf{x}}(t), \mathbf{0}) + \boldsymbol{\delta}(t),$$

$$(3.2) \quad \frac{D\boldsymbol{\delta}}{Dt} = \mathbf{G}(\hat{\mathbf{x}}(t-s_0), \dots, \hat{\mathbf{x}}(t-s_p), \boldsymbol{\delta}(t-s_0), \dots, \boldsymbol{\delta}(t-s_p)),$$

where $\boldsymbol{\delta} \in \mathbb{R}^Q$ is the closure term, $p \in \mathbb{N}$ is the number of delays of past time information, and $\frac{D}{Dt}$ represents time discretization. $\{s_i\}_{i=0}^p$ is given as a strictly monotonic equally spaced time sequence with $s_i = i\Delta t$, $s_i \in [0, t]$.

We note that Shulkind et al. [41] also pursue closure, but are focused on developing a Markovian correction term with the restriction that the magnitude of the closure term is small compared to the resolved term. In the present work, memory effects are represented via an additional governing equation for $\boldsymbol{\delta}$. Two important features of the current framework are: a) the functional form of \mathbf{G} is extracted from data, i.e., solution snapshots from the

high-dimensional model, and b) this framework is inherently non-Markovian (for the resolved variables $\hat{\mathbf{x}}$). As a side note in [Appendix A](#), the current framework of operator inference for the shortest memory case $p = 0$ is compared to Elman's network [16], a forerunner to modern recurrent neural networks [19], to highlight similarities and differences.

3.1. Interpretation of operator inference. In a simple setting, assume the dynamical system above is discretized using first-order forward time integration. Rewrite the partitioned system as

$$(3.3) \quad \frac{\hat{\mathbf{x}}^{n+1} - \hat{\mathbf{x}}^n}{\Delta t} = \hat{\mathbf{F}}(\hat{\mathbf{x}}^n, \mathbf{0}) + \hat{\mathbf{F}}(\hat{\mathbf{x}}^n, \tilde{\mathbf{x}}^n) - \hat{\mathbf{F}}(\hat{\mathbf{x}}^n, \mathbf{0}) = \hat{\mathbf{F}}(\hat{\mathbf{x}}^n, \mathbf{0}) + \delta^n,$$

$$(3.4) \quad \frac{\tilde{\mathbf{x}}^{n+1} - \tilde{\mathbf{x}}^n}{\Delta t} = \tilde{\mathbf{F}}(\hat{\mathbf{x}}^n, \tilde{\mathbf{x}}^n),$$

where $\delta = \hat{\mathbf{F}}(\hat{\mathbf{x}}, \tilde{\mathbf{x}}) - \hat{\mathbf{F}}(\hat{\mathbf{x}}, \mathbf{0}) \triangleq \mathbf{R}(\hat{\mathbf{x}}, \tilde{\mathbf{x}})$. Note that $\delta^{n+1} = \hat{\mathbf{F}}(\hat{\mathbf{x}}^{n+1}, \tilde{\mathbf{x}}^{n+1}) - \hat{\mathbf{F}}(\hat{\mathbf{x}}^{n+1}, \mathbf{0}) = \mathbf{R}(\hat{\mathbf{x}}^{n+1}, \tilde{\mathbf{x}}^{n+1}) = \mathbf{R}(\hat{\mathbf{x}}^n + \Delta t \hat{\mathbf{F}}(\hat{\mathbf{x}}^n, \mathbf{0}) + \Delta t \delta^n, \tilde{\mathbf{x}}^n + \Delta t \tilde{\mathbf{F}}(\hat{\mathbf{x}}^n, \tilde{\mathbf{x}}^n))$. Thus, one must obtain $\tilde{\mathbf{x}}^n$ to further evolve the closure.

As implied by the Takens embedding theorem, it is possible to use the information of past resolved states to obtain $\tilde{\mathbf{x}}^n$. Considering a time delay up to p steps, the equations that involve $\tilde{\mathbf{x}}$ are given as follows:

$$(3.5) \quad \delta^n = \mathbf{R}(\hat{\mathbf{x}}^n, \tilde{\mathbf{x}}^n),$$

$$(3.6) \quad \frac{\tilde{\mathbf{x}}^n - \tilde{\mathbf{x}}^{n-1}}{\Delta t} = \tilde{\mathbf{F}}(\hat{\mathbf{x}}^{n-1}, \tilde{\mathbf{x}}^{n-1}),$$

$$(3.7) \quad \delta^{n-1} = \mathbf{R}(\hat{\mathbf{x}}^{n-1}, \tilde{\mathbf{x}}^{n-1}),$$

⋮

$$(3.8) \quad \frac{\tilde{\mathbf{x}}^{n-p+1} - \tilde{\mathbf{x}}^{n-p}}{\Delta t} = \tilde{\mathbf{F}}(\hat{\mathbf{x}}^{n-p}, \tilde{\mathbf{x}}^{n-p}),$$

$$(3.9) \quad \delta^{n-p} = \mathbf{R}(\hat{\mathbf{x}}^{n-p}, \tilde{\mathbf{x}}^{n-p}),$$

with the number of equations, component-wise, being $N_{eq} = pN + Q$ and the number of unknowns, component-wise, being $N_{unk} = (p+1)(N-Q)$. Note that $N_{eq} - N_{unk} = (p+2)Q - N$. Therefore, for large enough p , it should be possible to determine $\tilde{\mathbf{x}}^n$ from $\hat{\mathbf{x}}^{n-1}, \dots, \hat{\mathbf{x}}^{n-p}$ and $\delta^n, \dots, \delta^{n-p}$ by solving the nonlinear algebraic equations above. Once $\tilde{\mathbf{x}}^n$ is determined, one can obtain $\frac{\delta^{n+1} - \delta^n}{\Delta t}$ as a function \mathbf{G} of $\hat{\mathbf{x}}^n, \dots, \hat{\mathbf{x}}^{n-p}$ and $\delta^n, \dots, \delta^{n-p}$. This suggests the possibility of finding \mathbf{G} through a data-driven method.

3.2. Definition and data preparation. As discussed above, the goal of the operator inference framework is to determine \mathbf{G} in (3.2). This process can also be viewed as a nonlinear system identification problem by considering δ as the states of an undetermined system and $\hat{\mathbf{x}}$ as inputs to this system. Our approach requires the parameterization of $\mathbf{G}(\cdot)$ in the form of $\mathbf{G}_{\mathbf{W}}$ using two different methodologies and then solving an optimization problem over the parameter space \mathbf{W} . The first parametrization method is sparse polynomial regression (similar to the SINDy approach of Brunton et al. [8]) which leverages the fact that many

dynamical systems in science and engineering can be represented as a sparse combination of monomials. The second method uses time delayed neural networks [51] which are scalable to high-dimensional nonlinear systems and possesses the universal approximator property and implicit feature selection. Note that for simplicity, the time derivative is realized by a first-order forward scheme throughout this work.

Assume M temporally sequential snapshots $\mathbf{X} = [\hat{\mathbf{x}}^j]_{j \in I} \in \mathbb{R}^{M \times Q}$ spaced uniformly with a time interval Δt , i.e., $s_i = i\Delta t$, $\forall i \in \{0, \dots, p\}$ and $d\mathbf{X} = [D\hat{\mathbf{x}}^j/Dt]_{j \in I} \in \mathbb{R}^{M \times Q}$ obtained from the full order model (Eq.(2.1)). Here, $I = \{j | j \in \mathbb{N}^+, 1 \leq j \leq M\}$, $M \in \mathbb{N}^+$ and $p \in \mathbb{N}$ is the number of steps of past memory. We divide \mathbf{X} into training and testing data through the index set, considering p time delays: $I^p = \{j | j \in I, \forall i \in \mathbb{N}, 0 \leq i \leq p, j - i \in I\}$; training data index set: $I_{train}^p = \{j | j \in I^p, j \leq M_{train}\}$, $M_{train} \in \mathbb{N}^+$; testing data index set: $I_{test}^p = I^p \setminus I_{train}^p$. It should be noted that we simply choose $s_i = i\Delta t$, where Δt is the time interval between equally spaced snapshots. As a result, given $\hat{\mathbf{F}}(\cdot, \mathbf{0})$, the corresponding snapshots of training closure are

$$(3.10) \quad \Delta_{I_{train}^p} = [\boldsymbol{\delta}^j]_{j \in I_{train}^p} = [D\hat{\mathbf{x}}^j/Dt - \hat{\mathbf{F}}(\hat{\mathbf{x}}^j, \mathbf{0})]_{j \in I_{train}^p}.$$

Therefore, the time-delayed feature matrix of $\hat{\mathbf{x}}$ and $\boldsymbol{\delta}$ in the training data can be constructed as

$$(3.11) \quad \mathbf{Y}_{I_{train}^p} = [\hat{\mathbf{x}}^j, \dots, \hat{\mathbf{x}}^{j-p}, \boldsymbol{\delta}^j, \dots, \boldsymbol{\delta}^{j-p}]_{j \in I_{train}^p} = [\mathbf{y}^j]_{j \in I_{train}^p},$$

where $\mathbf{y}^j \in \mathbb{R}^{2(1+p)Q}$.

Considering the dependency indicated by the relation between the sequences of $\boldsymbol{\delta}$ and $\hat{\mathbf{x}}$, $\forall j \in \mathbb{N}^+$, $j \leq n$, $\boldsymbol{\delta}^{n-j} = \frac{\hat{\mathbf{x}}^{n-j+1} - \hat{\mathbf{x}}^{n-j}}{\Delta t} - \hat{\mathbf{F}}(\hat{\mathbf{x}}^{n-j}, \mathbf{0})$; the economic time-delayed feature matrix can be constructed as

$$(3.12) \quad \mathbf{Y}_{I_{train}^p}^{eco} = [\hat{\mathbf{x}}^j, \dots, \hat{\mathbf{x}}^{j-p}, \boldsymbol{\delta}^j]_{j \in I_{train}^p} = [\mathbf{y}_{eco}^j]_{j \in I_{train}^p},$$

where $\mathbf{y}_{eco}^j \in \mathbb{R}^{(2+p)Q}$.

The training target is

$$(3.13) \quad \mathbf{Z}_{I_{train}^p} = D\Delta_{I_{train}^p}/Dt = [D\boldsymbol{\delta}^j/Dt]_{j \in I_{train}^p} = [\mathbf{z}^j]_{j \in I_{train}^p},$$

where $\mathbf{z}^j \in \mathbb{R}^Q$.

Likewise, the testing feature matrix and target are

$$(3.14) \quad \mathbf{Y}_{I_{test}^p} = [\hat{\mathbf{x}}^j, \dots, \hat{\mathbf{x}}^{j-p}, \boldsymbol{\delta}^j, \dots, \boldsymbol{\delta}^{j-p}]_{j \in I_{test}^p} = [\mathbf{y}^j]_{j \in I_{test}^p},$$

$$(3.15) \quad \mathbf{Z}_{I_{test}^p} = \frac{d\Delta}{dt} = [\mathbf{z}^j]_{j \in I_{test}^p},$$

and the corresponding economic feature matrix is

$$(3.16) \quad \mathbf{Y}_{I_{test}^p}^{eco} = [\hat{\mathbf{x}}^j, \dots, \hat{\mathbf{x}}^{j-p}, \boldsymbol{\delta}^j]_{j \in I_{test}^p} = [\mathbf{y}_{eco}^j]_{j \in I_{test}^p}.$$

3.3. Data driven modeling. Based on the above definitions, the general idea is to transform the functional approximation problem into an optimization problem either through sparse polynomial regression (SINDy) or neural networks as described in the following subsections.

3.3.1. Sparse polynomial model. Since polynomial features frequently appear in many science and engineering applications, polynomial regression [8] is typically a popular choice. An approximation of the form $\mathbf{G} = \hat{\mathbf{G}}\mathbf{W} = \Theta_k\mathbf{W}$ is employed by transforming the problem of finding a sparse representation of dynamical system into a convex optimization problem in (3.18), where $\mathbf{W} \in \mathbb{R}^{L_k^p \times Q}$ is a matrix of decision variables, and $L_k^p \in \mathbb{N}^+$ is the number of component-wise polynomial features up to total degree k of resolved states and using previous p steps. To illustrate this idea, given row vector $\mathbf{h} \in \mathbb{R}^{1 \times n}$, $\forall n \in \mathbb{N}^+$, Θ_k is a corresponding feature row vector from a monomial library with a certain maximum total order k

$$(3.17) \quad \Theta_k(\mathbf{h}) = [1 \quad \mathbf{h} \quad \mathbf{h}^{P_2} \dots \mathbf{h}^{P_k}],$$

where \mathbf{h}^{P_k} refers to all product terms of monomials with total degree k . Naturally, $\forall m, n \in \mathbb{N}^+$, $\mathbf{H} \in \mathbb{R}^{m \times n}$, $\Theta_k(\mathbf{H})$ is the row stack of $\Theta_k(\mathbf{H}_i)$, $i \in \{1, \dots, m\}$, $\mathbf{H} = [\mathbf{H}_1^T, \dots, \mathbf{H}_m^T]^T$.

The basic idea of sparse polynomial regression is to find a sparse representation by sparsifying the coefficient matrix \mathbf{W} from a predefined feature library through either a sequential thresholded least-squares algorithm [8] or using *lasso* [29].

Applying *lasso* [46, 47], given $|I_{train}^p| \geq (1+p)L_k^p$, (3.18) can be shown to be a convex optimization problem for \mathbf{w}_k , the k th column of \mathbf{W} :

$$(3.18) \quad \mathbf{w}_k^* = \underset{\mathbf{w}_k}{\operatorname{argmin}} \frac{1}{|I_{train}^p|} \sum_{j \in I_{train}^p} \|\mathbf{z}^j - \Theta_k(\mathbf{y}^j)\mathbf{w}_k\|_2^2 + \lambda \|\mathbf{w}_k\|_1,$$

where $|\cdot|$ is the cardinality, $\Theta_k(\mathbf{y}^j) \in \mathbb{R}^{|I_{train}^p| \times L_k^p}$, and $\lambda \in \mathbb{R}^+$ is the penalty coefficient. It is important to note that for *lasso* to achieve ideal support recovery, the following constraint must be satisfied [44]

$$(3.19) \quad n_k/n_p \leq \delta/(2 \log n_p)(1 + o(1)),$$

where n_k is the number of true features, n_p is the number of total features, $\delta = n/n_p$ where n is the number of i.i.d data points. The simplest way to achieve this is to ensure one has a large number data n compared to number of total features n_p . *lasso* is implemented via Scikit-learn packages [34].

Regarding the dependency between features, since $j \leq n$, $\delta^{n-j} = \frac{\hat{\mathbf{x}}^{n-j+1} - \hat{\mathbf{x}}^{n-j}}{\Delta t} - \hat{\mathbf{F}}(\hat{\mathbf{x}}^{n-j}, \mathbf{0})$. Hence, if $\hat{\mathbf{F}}$ is compact in monomials, it is possible to replace $\delta^{n-1}, \dots, \delta^{n-p}$ by polynomials of $\hat{\mathbf{x}}^n, \dots, \hat{\mathbf{x}}^{n-p}$. However, if $\hat{\mathbf{F}}$ is not in polynomial form or if $\hat{\mathbf{F}}$ contains a very high order polynomial of $\hat{\mathbf{x}}$, the size of the library will be extremely large and potentially non-convergent. While reuse of δ^{n-j} can alleviate this issue, the trade-off involves using twice the degrees of freedom. For the sparse polynomial regression model, employment of δ^{n-j} is considered throughout this work for better fitting performance.

The current framework with polynomial features is different from the operator inference of Peherstorfer and Willcox [35], which targets the entire system and can be viewed as a particular

case of the present work with $p = 0$, $k = 2$, and $\hat{\mathbf{F}} = 0$, and with POD as preprocessing for dimension reduction. Additionally, sparsity is not encouraged and no memory effects are required in their model. The present framework thus seeks a more general non-Markovian operator inference.

3.3.2. Artificial neural network model. In the previous subsection, the number of polynomial features in the feature library is found to grow exponentially with Q . One of the most popular tools for efficient high-dimensional functional approximations is the artificial neural network (ANN). The appealing feature of neural networks with a single hidden layer and squashing nonlinear activation function is that any Borel measurable function can be approximated to any degree of accuracy on a compact domain. This property is guaranteed by the universal approximation theory [23]. ANN has attracted considerable attention in recent years. The success of deep learning (naïvely and narrowly speaking for supervised learning, ANN with a large number of hidden layers) lies in learning low-dimensional representations from high-dimensional, complex data effectively and building relationships between learned features and the target [19].

To parametrize the model described in (3.24), the standard feedforward neural network structure shown in Figure 1 is employed for $\mathbf{G} = \hat{\mathbf{G}}_{\boldsymbol{\theta}, \mathbf{b}}$. Due to the previously mentioned dependency between sequences of $\boldsymbol{\delta}$ and $\hat{\mathbf{x}}$ and the universal approximator property of ANN, $\mathbf{y}_{eco} \in \mathbb{R}^{(2+p)Q}$ is used as input. To construct a densely connected feedforward neural network $\hat{\mathbf{G}}_{\boldsymbol{\theta}, \mathbf{b}}: \mathbb{R}^{(2+p)Q} \mapsto \mathbb{R}^Q$ with $L-1$ hidden layers and a linear output layer, the following recursive expression is used for each hidden layer:

$$(3.20) \quad \boldsymbol{\eta}^l = \sigma_l(\theta_l \boldsymbol{\eta}^{l-1} + b_l),$$

for $l = 1, \dots, L-1$, where $\boldsymbol{\eta}^0$ stands for the input of the neural network, $\boldsymbol{\eta}^l \in \mathbb{R}^{n_l \times 1}$, $n_l \in \mathbb{N}^+$ is the number of units in layer l , $\theta_l \in \mathbb{R}^{n_l \times n_{l-1}}$, $b_l \in \mathbb{R}^{n_l \times 1}$, $n_0 = (2+p)Q$, and σ_l is the activation function of layer l . Note that the output layer is linear, i.e., $\sigma_L(x) = x$:

$$(3.21) \quad \hat{\mathbf{G}}(\mathbf{y}_{eco}; \boldsymbol{\theta}, \mathbf{b}) = \boldsymbol{\eta}^L = \theta_L \boldsymbol{\eta}^{L-1} + b_L,$$

where $\theta_L \in \mathbb{R}^{n_L \times n_{L-1}}$, $b_L \in \mathbb{R}^{n_L \times 1}$, and $n_L = Q$. Parameters of the neural network are summarized as $\mathbf{W} = \{\boldsymbol{\theta}, \mathbf{b}\}$: weights $\boldsymbol{\theta} = \{\theta_j\}_{j=1, \dots, L}$ and biases $\mathbf{b} = \{b_j\}_{j=1, \dots, L}$. In this work we use two hidden layers where $L = 3$ and hidden units are all the same. The full expression of the neural network model is

$$(3.22) \quad \hat{\mathbf{G}}(\mathbf{y}_{eco}; \mathbf{W}) = \hat{\mathbf{G}}(\mathbf{y}_{eco}; \boldsymbol{\theta}, \mathbf{b}) = \theta_3 \sigma(\theta_2 \sigma(\theta_1 \mathbf{y}_{eco} + b_1) + b_2) + b_3,$$

where $\sigma(\cdot): \mathbb{R} \mapsto \mathbb{R}$ is a nonlinear activation function, e.g., ReLU, SeLU, tanh [1].

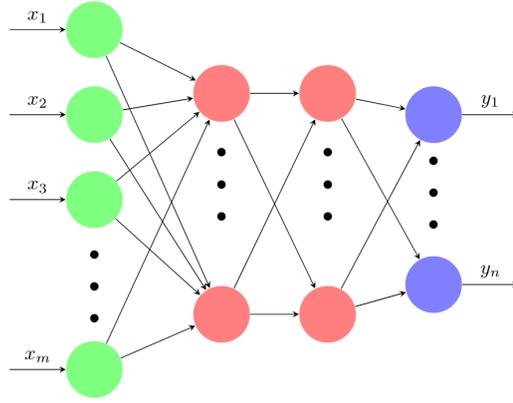


Figure 1. Schema of a typical feedforward neural network $\mathbb{R}^m \mapsto \mathbb{R}^n$ with two hidden layers with \mathbf{x} as input and \mathbf{y} as output

The problem of finding a good neural network model is equivalent to searching for a set of parameters $\boldsymbol{\theta}$ and \mathbf{b} that optimize the mean-square-error on training data with weight decay regularization

$$(3.23) \quad \mathbf{W}^* = \{\boldsymbol{\theta}^*, \mathbf{b}^*\} = \underset{\boldsymbol{\theta}, \mathbf{b}}{\operatorname{argmin}} \frac{1}{|I_{train}^p|} \sum_{j \in I_{train}^p} \left\| \mathbf{z}^j - \hat{\mathbf{G}}(\mathbf{y}_{eco}^j; \boldsymbol{\theta}, \mathbf{b}) \right\|_2^2 + \lambda \sum_{l=1}^L \|\boldsymbol{\theta}_l\|_F^2,$$

where $\|(\cdot)\|_F$ denotes the Frobenius norm. The weights are initialized using the standard truncated normal distribution, and a first order gradient-based technique [26] is used for optimization. Unfortunately, due to the non-convex nature of (3.22), one can often only afford to find a local minimum instead of the global minimum. However, in practice, a local minimum is usually satisfactory if it is properly trained and validated. Hyperparameters for each case given below are selected using grid search in a certain range. The model is implemented with the Keras [11] and Tensorflow libraries [1].

3.4. Reducing the computational complexity of multi-time effects. From a training perspective, the most difficult part of generating the polynomial model is in characterizing multi-time effects. The most general way of treating the memory effect in (3.1) is to consider all interactions between the past, i.e., cross-time memory effects, similar to a nonlinear autoregressive network with exogenous inputs (NARX) model [7] with multi-time correlations. Unfortunately, this method of discovering \mathbf{G} is severely hindered by the curse of dimensionality. Introducing full interactions with polynomial features up to a total degree $k \in \mathbb{N}^+$ would lead to a number of features scaling as $L_k^p = \binom{2(1+p)Q+k}{k} \propto (2(1+p)Q)^k$.

An alternative strategy to build computationally feasible set of features is to *assume* that full memory effects can be approximated with linear superposition of multi-time nonlinear features, i.e., linear treatment of multi-time memory in the form

$$(3.24) \quad \mathbf{G}(\hat{\mathbf{x}}(t-s_0), \dots, \hat{\mathbf{x}}(t-s_p), \boldsymbol{\delta}(t-s_0), \dots, \boldsymbol{\delta}(t-s_p); \mathbf{W}) \approx \sum_{i=0}^p \mathbf{G}_i(\hat{\mathbf{x}}(t-s_i), \boldsymbol{\delta}(t-s_i); \mathbf{W}^i),$$

where $\mathbf{G}_i(\cdot, \cdot) : \mathbb{R}^Q \times \mathbb{R}^Q \mapsto \mathbb{R}^Q$ represents the contribution of the system at $t = t - s_i$ to t in closure dynamics. With a polynomial basis, this approach will reduce the number of features up to k total degrees to $L_{k, reduced}^p = (1 + p) \binom{2Q+k}{k} - p \propto p(2Q)^k$, which grows linearly with p ¹.

The above assumption would lead to a reduction in the number of fitting parameters with respect to increasing memory length p . We note that, if this decoupling is applied to the time delay neural network (TDNN) model, this strategy can be viewed as a regularization of the neural network model by pruning weights between units of different time instances.

3.5. Model selection. Since both types of models mentioned above require the specification of hyperparameters before training, model selection is an important aspect. For the sparse polynomial regression model, the associated hyperparameters are:

- maximum total degree of polynomials: k
- maximum number of previous states: p
- penalty coefficient: λ

In practice, we choose p and k heuristically and as small as possible, while still fitting the model with ordinary-least-square (OLS) and keeping the total number of features smaller than the number of samples to ensure strict convexity of the *lasso* problem [47]. To determine λ , we draw the *lasso* path to decide the most appropriate solution that balances complexity and mean-squared-error (MSE) error. It should be noted that when drawing the *lasso* path, we split the training data in time; the first 80% is training data used to compute the *lasso* path, and the last 20% is validation data. The goal is to obtain a robust model that generalizes beyond the training set.

For the neural network, a logical strategy of hyperparameter selection has proved to be challenging for even the simplest standard feedforward neural network. In the present work, since the problem size is small, we choose hyperparameters via simple grid search for the type of activation function and number of hidden units.

3.6. Evaluation of MSE a priori and a posteriori. Notice that the optimization problems described in (3.18) and (3.23) only guarantee performance in an a priori sense on training data. A proper evaluation of the model should be performed both in an a priori sense as mean-squared-error over the data index set I^p ,

$$(3.25) \quad e_{\text{MSE}}^{\text{apr}} = \frac{1}{|I^p|} \sum_{j \in I^p} \left\| \mathbf{z}^j - \hat{\mathbf{G}}(\mathbf{y}^j; \mathbf{W}^*) \right\|_2^2,$$

and in an a posteriori sense in which only the initial condition is given to the model, as

$$(3.26) \quad e_{\text{MSE}}^{\text{apo}} = \frac{1}{|I^p|} \sum_{j \in I^p} \left\| \hat{\mathbf{x}}^j - \hat{\mathbf{x}}^{*j} \right\|_2^2,$$

where $\hat{\mathbf{x}}^*(t)$ is the solution of the augmented system (3.1) and (3.2) with $\hat{\mathbf{G}}(\cdot; \mathbf{W}^*)$ starting with an exact initial condition. This type of a posteriori evaluation is also called free-run in the time series modeling community [2].

¹– p comes from removing the redundant constant feature

4. Results - Linear system. To illustrate the idea of applying operator inference and to motivate further developments, the polynomial closure model is first applied on a three dimensional linear system shown below

$$(4.1) \quad \frac{d}{dt} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 & -1 & -1 \\ 0.5 & -1.1 & 1.5 \\ 1 & -3 & 0.5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix},$$

where $\hat{x}(0) = \hat{x}^0 = x_1^0 = 1$ and $x_2^0 = x_3^0 = 0$. The first-order forward discretized form of (4.1) is shown in (4.2) with total degrees of freedom $N = 3$ and number of reduced states $Q = 1$

$$(4.2) \quad \begin{bmatrix} x_1^{n+1} \\ x_2^{n+1} \\ x_3^{n+1} \end{bmatrix} = \begin{bmatrix} x_1^n \\ x_2^n \\ x_3^n \end{bmatrix} + \Delta t \begin{bmatrix} 0 & -1 & -1 \\ 0.5 & -1.1 & 1.5 \\ 1 & -3 & 0.5 \end{bmatrix} \begin{bmatrix} x_1^n \\ x_2^n \\ x_3^n \end{bmatrix}.$$

Consequently, following the operator inference framework with the polynomial form in (3.1) and (3.2) and a linear superposition assumption of multi-time effects, we have $\hat{x} = x_1$, $\hat{F} = 0$ for the following ROM formulation

$$(4.3) \quad \hat{x}^{n+1} = \hat{x}^n + \Delta t \delta^n,$$

$$(4.4) \quad \delta^{n+1} = \delta^n + \Delta t \sum_{i=0}^p \mathbf{G}_i.$$

The goal is to extract a functional form of the governing equation $\{\mathbf{G}_i\}_{i=0}^p$ for closure δ from data, i.e., to determine $(\delta^{n+1} - \delta^n)/\Delta t$ as a function of previous \hat{x} and δ from data. The true closure is $\delta = -x_2 - x_3$, $\tilde{\mathbf{x}} = [x_2^\top, x_3^\top]^\top$, which is assumed to be unknown to the ROM. The simulation is run for $t \in [0, 40]$ and $\Delta t = 0.01$, resulting in a collection of 4000 snapshots in $\{\hat{x}, \delta\}$. The first 10% of data is set for training and the remaining 90% as testing data.

For this 3D linear system, the exact solution for the closure dynamics is

$$(4.5) \quad \frac{\delta^{n+1} - \delta^n}{\Delta t} = \left(\frac{3}{2} - \frac{17\Delta t}{20} \right) \hat{x}^{n-1} - \left(\frac{183\Delta t + (35\Delta t + 10) \left(\frac{1}{\Delta t} - \frac{41}{10} \right)}{10} \right) \delta^{n-1} - \frac{3}{2} \hat{x}^n + \left(\frac{35\Delta t + 10}{10\Delta t} - \frac{41}{10} \right) \delta^n.$$

4.1. Model selection. As displayed in Figures 2 and 3, by applying the model with $p = 1$, $k = 1$ and λ chosen as 10^{-12} from the Pareto front of *lasso* path, the resulting model is found to only contain 4 non-zero terms.

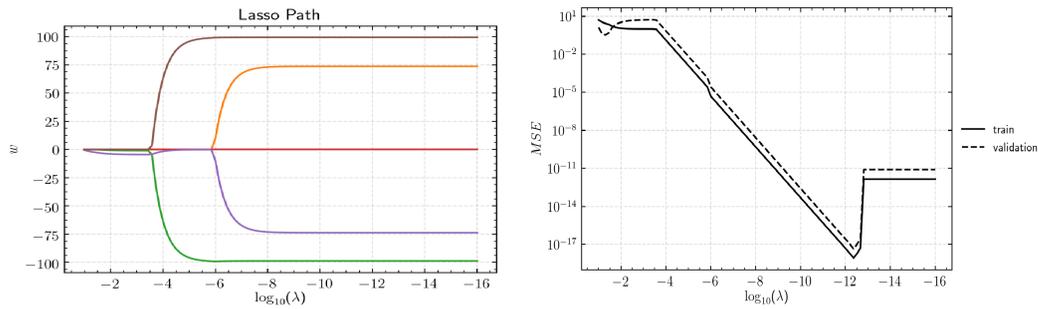


Figure 2. lasso path for the 3D linear system. Left: coefficients. Right: MSE.

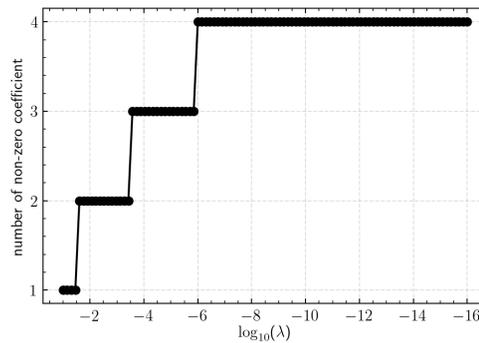


Figure 3. lasso path for the 3D linear system: number of non-zero terms

4.2. A posteriori evaluation of model performance. Using the hyperparameters determined above, the predicted trajectory of $\hat{x}(t)$ is found to match the target trajectory to an excellent degree, as shown in Figure 4.

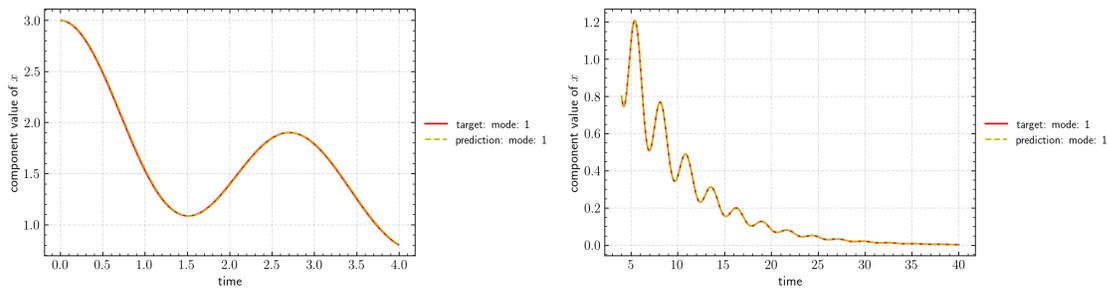


Figure 4. A posteriori model performance on the linear system. Left: training data. Right: testing data.

However, if one sets $p = 0$, the resulting model cannot produce good predictions, as the true solution is $\delta^{n+1} - \delta^n = -\Delta t(1.5\hat{x}^n + 4.1\delta^n + 6.1x_3^n)$. Since x_3^n is unknown to $\{\hat{x}^n, \delta^n\}$, additional memory length is required. The lasso path is shown in Figures 5 and 6 for the case with insufficient memory. The corresponding a posteriori performance is shown to be poor in Figure 7.

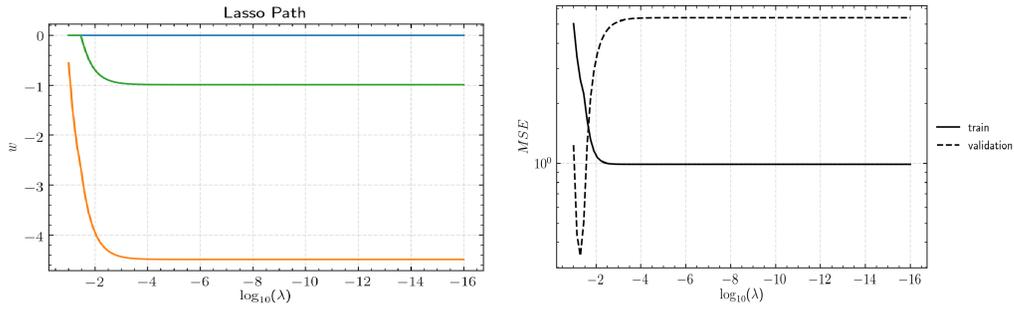


Figure 5. lasso path for the 3D linear system. Left: coefficients. Right: MSE.

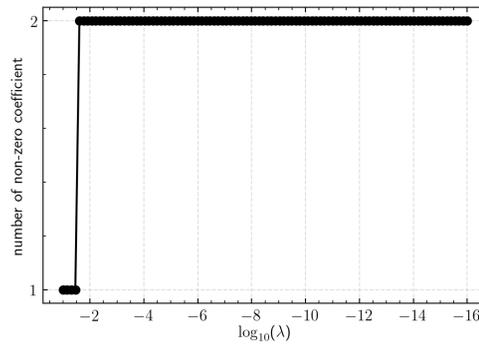


Figure 6. lasso path for the 3D linear system: number of non-zero terms

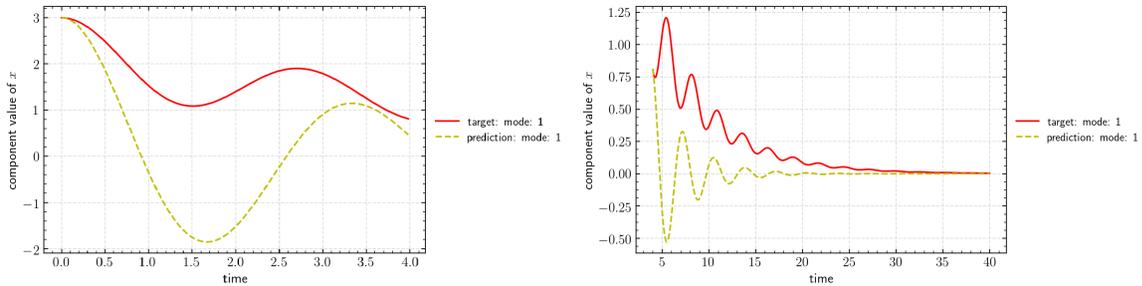


Figure 7. A posteriori model performance on the linear system with $p = 0$. Left: training data. Right: testing data.

One might suspect that the sparsest solution, i.e., $\min \|\mathbf{w}_k\|_0$, should contain at most 3 non-zero terms because as mentioned before, there is one redundancy in $\{x_1^n, x_1^{n-1}, \delta^{n-1}\}$. The lasso based on the ℓ_1 norm does not guarantee the sparsest solution in the sense of the ℓ_0 norm, but it makes the problem computationally tractable.

5. Theoretical results. In this section, theoretical results are presented with regard to the capability of the closure model to determine the resolved dynamics with time-delayed features.

Definition 5.1 (Nonlinear dynamical system with dual linear closure). A nonlinear dynamical system with dual linear closure is defined in the following form:

$$\frac{d}{dt} \begin{bmatrix} \hat{\mathbf{x}} \\ \tilde{\mathbf{x}} \end{bmatrix} = \begin{bmatrix} \mathbf{F}(\hat{\mathbf{x}}) + A_{12}\tilde{\mathbf{x}} \\ \mathbf{H}(\hat{\mathbf{x}}) + A_{22}\tilde{\mathbf{x}} \end{bmatrix},$$

where $\hat{\mathbf{x}} \in \mathbb{R}^Q$, $\tilde{\mathbf{x}} \in \mathbb{R}^{N-Q}$, $\mathbf{F}(\cdot) : \mathbb{R}^Q \mapsto \mathbb{R}^Q$ and $\mathbf{H}(\cdot) : \mathbb{R}^Q \mapsto \mathbb{R}^{N-Q}$, $A_{12} \in \mathbb{R}^{Q \times N-Q}$, $A_{22} \in \mathbb{R}^{N-Q \times N-Q}$ with $\boldsymbol{\delta} = A_{12}\tilde{\mathbf{x}}$, $N \in \mathbb{N}$ and $Q \in \mathbb{N}$, $Q < N$.

The corresponding first order forward discretized dynamical system is

$$\begin{bmatrix} \hat{\mathbf{x}}^{n+1} \\ \tilde{\mathbf{x}}^{n+1} \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{x}}^n \\ \tilde{\mathbf{x}}^n \end{bmatrix} + \Delta t \begin{bmatrix} \mathbf{F}(\hat{\mathbf{x}}^n) + A_{12}\tilde{\mathbf{x}}^n \\ \mathbf{H}(\hat{\mathbf{x}}^n) + A_{22}\tilde{\mathbf{x}}^n \end{bmatrix},$$

where n denotes steps in time.

The exact discrete closure dynamics is

$$\boldsymbol{\delta}^{n+1} = A_{12}\tilde{\mathbf{x}}^{n+1} = \boldsymbol{\delta}^n + \Delta t(A_{12}A_{22}\tilde{\mathbf{x}}^n + A_{12}\mathbf{H}(\hat{\mathbf{x}}^n)).$$

Note that the only unknown is $\tilde{\mathbf{x}}^n$. Clearly, a sufficient condition for the data-driven framework to exactly represent the closure term $A_{12}\tilde{\mathbf{x}}^n$, would be the determination of $\tilde{\mathbf{x}}^n$ from a linear combination of $\hat{\mathbf{x}}^{n-1}, \dots, \hat{\mathbf{x}}^{n-p}$, $\boldsymbol{\delta}^n, \dots, \boldsymbol{\delta}^{n-p}$. It will be shown that, with certain structures of A_{12} and A_{22} , one can recover the entire history of $\tilde{\mathbf{x}}$ using up to previous p steps. The following proposition has strong similarities with the observability problem in linear system theory.

Theorem 5.2. For $k \in \mathbb{N}^+$, define the following matrix $\mathcal{O}_k \in \mathbb{R}^{kQ \times (N-Q)}$

$$\mathcal{O}_k = \begin{bmatrix} A_{12} \\ A_{12}A_{22} \\ \vdots \\ A_{12}A_{22}^{k-1} \end{bmatrix},$$

and the following mapping $r_{\mathcal{O}}(\cdot) : \mathbb{N}^+ \mapsto \mathbb{N}$

$$r_{\mathcal{O}}(k) = \text{rank}(\mathcal{O}_k).$$

If \mathcal{O}_{N-Q} is full rank, i.e., $r_{\mathcal{O}}(N-Q) = N-Q$, then for the first order forward discretized system with dual linear closure, $\exists p, n \in \mathbb{N}$, with collected $\boldsymbol{\delta} \in \mathbb{R}^Q$ and $\hat{\mathbf{x}} \in \mathbb{R}^Q$ up to step n , such that $\tilde{\mathbf{x}}^n, \dots, \tilde{\mathbf{x}}^{n-p}$ can be determined as a linear combination of $\mathbf{H}(\hat{\mathbf{x}}^{n-1}), \dots, \mathbf{H}(\hat{\mathbf{x}}^{n-p})$, $\boldsymbol{\delta}^n, \dots, \boldsymbol{\delta}^{n-p}$. For $p=0$, only $\boldsymbol{\delta}^n$ is used.

Further, the minimal p_* that satisfies the above is

$$p_* = \min \Omega_{\mathcal{O}} - 1,$$

where

$$\Omega_{\mathcal{O}} = \{l \in \mathbb{N}^+, r_{\mathcal{O}}(l) = N-Q\}.$$

Proof. Consider the first order forward discretized system of a dynamical system with linear closure, with $n, p \in \mathbb{N}$, $n > p$. We have the following *evolution* equations for the unresolved variable $\tilde{\mathbf{x}}$

$$(5.1) \quad \tilde{\mathbf{x}}^n = (I + \Delta t A_{22}) \tilde{\mathbf{x}}^{n-1} + \Delta t \mathbf{H}(\hat{\mathbf{x}}^{n-1}),$$

...

$$(5.2) \quad \tilde{\mathbf{x}}^{n-p+1} = (I + \Delta t A_{22}) \tilde{\mathbf{x}}^{n-p} + \Delta t \mathbf{H}(\hat{\mathbf{x}}^{n-p}),$$

and *projection* equations for δ

$$(5.3) \quad A_{12} \tilde{\mathbf{x}}^n = \delta^n,$$

...

$$(5.4) \quad A_{12} \tilde{\mathbf{x}}^{n-p} = \delta^{n-p}.$$

Note that the independent unknowns are $\{\tilde{\mathbf{x}}^{n-i}\}_{i=0}^p$ and we are provided with $\{\hat{\mathbf{x}}^{n-i}\}_{i=1}^p$ and $\{\delta^{n-i}\}_{i=0}^p$. Rearranging equations in matrix form, we have

$$(5.5) \quad \mathbf{\Gamma}_p \mathbf{X}_p = \mathbf{\Sigma}_p,$$

where

$$(5.6) \quad \mathbf{\Gamma}_p = \begin{bmatrix} I & -(I + \Delta t A_{22}) & \dots & \mathbf{0} \\ \mathbf{0} & I & -(I + \Delta t A_{22}) & \dots \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \dots & I & -(I + \Delta t A_{22}) \\ A_{12} & \mathbf{0} & \mathbf{0} & \dots \\ \mathbf{0} & A_{12} & \mathbf{0} & \dots \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & A_{12} \end{bmatrix},$$

$$(5.7) \quad \mathbf{X}_p = \begin{bmatrix} \tilde{\mathbf{x}}^n \\ \tilde{\mathbf{x}}^{n-1} \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \tilde{\mathbf{x}}^{n-p} \end{bmatrix}, \quad \mathbf{\Sigma}_p = \begin{bmatrix} \Delta t \mathbf{H}(\hat{\mathbf{x}}^{n-1}) \\ \Delta t \mathbf{H}(\hat{\mathbf{x}}^{n-2}) \\ \vdots \\ \Delta t \mathbf{H}(\hat{\mathbf{x}}^{n-p}) \\ \vdots \\ \delta^n \\ \vdots \\ \delta^{n-p} \end{bmatrix},$$

Using row operations to remove the diagonal block matrix of A_{12} ,

$$(5.8) \quad \mathbf{\Gamma}_p \rightarrow \begin{bmatrix} I & -(I + \Delta t A_{22}) & \dots & \mathbf{0} \\ \mathbf{0} & I & -(I + \Delta t A_{22}) & \dots \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & I & -(I + \Delta t A_{22}) \\ \mathbf{0} & \mathbf{0} & \dots & A_{12}(I + \Delta t A_{22})^p \\ \mathbf{0} & \mathbf{0} & \dots & A_{12}(I + \Delta t A_{22})^{p-1} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & A_{12} \end{bmatrix},$$

Thus

$$(5.9) \quad \text{rank}(\mathbf{\Gamma}_p) = p(N - Q) + \text{rank} \begin{bmatrix} A_{12}(I + \Delta t A_{22})^p \\ A_{12}(I + \Delta t A_{22})^{p-1} \\ \vdots \\ A_{12} \end{bmatrix}.$$

Note that

$$(5.10) \quad \text{rank} \begin{bmatrix} A_{12}(I + \Delta t A_{22})^p \\ A_{12}(I + \Delta t A_{22})^{p-1} \\ \vdots \\ A_{12} \end{bmatrix} = \text{rank} \begin{bmatrix} A_{12} A_{22}^p \\ A_{12} A_{22}^{p-1} \\ \vdots \\ A_{12} \end{bmatrix} = \text{rank}(\mathcal{O}_{p+1}).$$

From basic linear algebra, it is known that $r_{\mathcal{O}}(\cdot)$ is bounded and monotonic, where $r_{\mathcal{O}}(\cdot) : \mathbb{N} \mapsto \mathbb{N}$. Also recall that \mathcal{O}_{N-Q} is full rank thus the following set Ω is not empty

$$(5.11) \quad \Omega_{\mathcal{O}} = \{l | l \in \mathbb{N}^+, r_{\mathcal{O}}(l) = N - Q\}.$$

Therefore setting

$$(5.12) \quad p_* = \min \Omega_{\mathcal{O}} - 1,$$

we will have

$$(5.13) \quad \text{rank}(\mathbf{\Gamma}_p) = p(N - Q) + (N - Q) = (p + 1)(N - Q),$$

indicating $\mathbf{\Gamma}_p$ is full column rank. Therefore, consider $\mathbf{\Gamma}_p^+$ as the *left* Moore-Penrose inverse of $\mathbf{\Gamma}_p$

$$(5.14) \quad \mathbf{\Gamma}_p^+ = (\mathbf{\Gamma}_p^\top \mathbf{\Gamma}_p)^{-1} \mathbf{\Gamma}_p,$$

and thus

$$(5.15) \quad \mathbf{X}_p = \mathbf{\Gamma}_p^+ \mathbf{\Gamma}_p \mathbf{X}_p = \mathbf{\Gamma}_p^+ \mathbf{\Sigma}_p. \quad \blacksquare$$

Again, note that once $\tilde{\mathbf{x}}^n$ is determined from past time instances of $\boldsymbol{\delta}$ and $\hat{\mathbf{x}}$ up to step p , the closure dynamics is fully determined as well. As an example, applying [Theorem 5.2](#) to the 3D discrete linear system described in (4.1), $A_{12} = \begin{bmatrix} -1 & -1 \end{bmatrix}$ and $A_{22} = \begin{bmatrix} -1.1 & 1.5 \\ -3 & 0.5 \end{bmatrix}$, $\text{rank}(A_{12}) = 1$, $\text{rank}\left(\begin{bmatrix} A_{12} \\ A_{12}A_{22} \end{bmatrix}\right) = \text{rank}\left(\begin{bmatrix} -1 & -1 \\ 4.1 & -2 \end{bmatrix}\right) = 2$. Therefore, $p_* = \min \Omega_{\mathcal{O}} - 1 = 1$. As a trivial observation, based on the [Theorem 5.2](#), one can immediately obtain the following proposition.

Proposition 5.3. *If A_{12} has full column rank, $\tilde{\mathbf{x}}^n$ can be determined with only $\boldsymbol{\delta}^n$.*

However, full observability on all past states of $\tilde{\mathbf{x}}$ is a very strong criterion to guarantee predictability of dual linear closure dynamics. Indeed, from a data-driven perspective, one only requires that the linear closure is in the p -time delayed observable space of $\tilde{\mathbf{x}}$ as defined in [Definition 5.4](#). Therefore, we now turn our focus to finding $A_{12}\tilde{\mathbf{x}}^n$ directly. First, a strict definition of p -time delayed linear observable space is given below.

Definition 5.4 (p -time delayed linear observable space). *For the first order forward discretized nonlinear dynamical system with dual linear closure, define the corresponding p time delayed linear observable space χ_p as*

$$\chi_p = \{\eta | \eta = v^\top \mathbf{X}_p, v \in \text{Im } V_p\},$$

where V_p is from the reduced singular value decomposition of $\mathbf{\Gamma}_p$

$$\mathbf{\Gamma}_p = U_p S_p V_p^\top,$$

with $U_p \in \mathbb{R}^{(pN+Q) \times r}$, $S_p \in \mathbb{R}^{r \times r}$, $V_p \in \mathbb{R}^{(p+1)(N-Q) \times r}$, $r = \text{rank}(\mathbf{\Gamma}_p)$.

Regarding the question of determining a general linear combination of \mathbf{X}_p from a p -time delayed observable space, the following lemma shows that a rank test can provide essential insight.

Lemma 5.5. *For a nonlinear dynamical system with dual linear closure, for any quantity $\xi \in \mathbb{R}^{q \times 1}$ that is a linear combination of \mathbf{X}_p characterized by C ,*

$$\xi = C^\top \mathbf{X}_p,$$

where $C \in \mathbb{R}^{(p+1)(N-Q) \times q}$, if

$$\text{rank}(V_p) = \text{rank}\left(\begin{bmatrix} C^\top \\ V_p^\top \end{bmatrix}\right),$$

then

$$\xi \in \chi_p,$$

i.e., ξ is observable with p -time delayed information of $\boldsymbol{\delta}$ and $\hat{\mathbf{x}}$.

Proof. $\because \text{rank}\left(\begin{bmatrix} C^\top \\ V_p^\top \end{bmatrix}\right) = \text{rank}([C \ V_p]) = \text{rank}(V_p^\top) = \text{rank}(V_p) \therefore C \subset \text{Im } V_p$
 $\therefore \xi = C^\top \mathbf{X}_p \in \chi_p$. ■

Given the above lemma, one can obtain a rank test criterion in [Theorem 5.6](#) for whether the closure dynamics of a nonlinear system with dual linear closure can be determined with p time delayed observable space. Furthermore, analysis of the rank of the augmented matrix provides further insights into the role of time delay in observation.

Theorem 5.6. *A nonlinear dynamical system with dual linear closure with $p = N - Q - 1$ will satisfy the following rank test*

$$\text{rank}(V_p) = \text{rank}\left(\begin{bmatrix} C^\top \\ V_p^\top \end{bmatrix}\right),$$

where $C^\top = [A_{12}A_{22} \ \mathbf{0}]$, and the closure dynamics is observable from p time delayed observable space, i.e., can be determined as a linear combination of $\mathbf{H}(\hat{\mathbf{x}}^{n-1}), \dots, \mathbf{H}(\hat{\mathbf{x}}^{n-p}), \boldsymbol{\delta}^n, \dots, \boldsymbol{\delta}^{n-p}$. Furthermore, the minimal number of previous steps p that satisfies the above condition can be found through

$$p_* = \min \Pi_{\mathcal{O}} - 1,$$

where

$$\Pi_{\mathcal{O}} = \{l \in \mathbb{N}^+, r_{\mathcal{O}}(l) = r_{\mathcal{O}}(l+1)\}.$$

Proof. To determine $\boldsymbol{\delta}^{n+1}$, $A_{12}A_{22}\tilde{\mathbf{x}}^n$ has to be in the p time delayed observable space. $\therefore \tilde{\mathbf{x}}^n = [A_{12}A_{22} \ \mathbf{0}] \mathbf{X}_p$. \therefore from [Lemma 5.5](#), if

$$\text{rank}(V_p) = \text{rank}\left(\begin{bmatrix} A_{12}A_{22} & \mathbf{0} \\ V_p^\top \end{bmatrix}\right),$$

then $\tilde{\mathbf{x}}^n$ is p time delayed linear observable. Since V_p^\top shares the same independent row space

as $\mathbf{\Gamma}$, augmenting $\mathbf{\Gamma}_p$ with C^\top will result in the same rank.

$$\begin{aligned}
\begin{bmatrix} \mathbf{\Gamma}_p \\ C^\top \end{bmatrix} &= \begin{bmatrix} I & -(I + \Delta t A_{22}) & \dots & \mathbf{0} \\ \mathbf{0} & I & -(I + \Delta t A_{22}) & \dots \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \dots & I & -(I + \Delta t A_{22}) \\ A_{12} & \mathbf{0} & \mathbf{0} & \dots \\ \mathbf{0} & A_{12} & \mathbf{0} & \dots \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & A_{12} \\ A_{12}A_{22} & \mathbf{0} & \dots & \mathbf{0} \end{bmatrix} \\
\rightarrow & \begin{bmatrix} I & -(I + \Delta t A_{22}) & \dots & \mathbf{0} \\ \mathbf{0} & I & -(I + \Delta t A_{22}) & \dots \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & I & -(I + \Delta t A_{22}) \\ \mathbf{0} & \mathbf{0} & \dots & A_{12}(I + \Delta t A_{22})^p \\ \mathbf{0} & \mathbf{0} & \dots & A_{12}(I + \Delta t A_{22})^{p-1} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & A_{12} \\ \mathbf{0} & \mathbf{0} & \dots & A_{12}A_{22}(I + \Delta t A_{22})^p \end{bmatrix} \\
\rightarrow \text{rank}\left(\begin{bmatrix} \mathbf{\Gamma}_p \\ C^\top \end{bmatrix}\right) &= p(N - Q) + \text{rank}\left(\begin{bmatrix} A_{12}A_{22}^p \\ A_{12}A_{22}^{p-1} \\ \vdots \\ A_{12} \\ A_{12}A_{22}^{p+1} \end{bmatrix}\right) = p(N - Q) + \text{rank}(\mathcal{O}_{p+2}) \\
\rightarrow \text{rank}\left(\begin{bmatrix} \mathbf{\Gamma}_p \\ C^\top \end{bmatrix}\right) &= \text{rank}(\mathbf{\Gamma}_p) \rightarrow \text{rank}(\mathcal{O}_{p+2}) = \text{rank}(\mathcal{O}_{p+1}).
\end{aligned}$$

\therefore Recall $r_{\mathcal{O}}(\cdot)$ is a monotonic integer function and from the Cayley Hamilton theorem, A_{22}^{N-Q} is linearly dependent on $\{I, A_{22}, \dots, A_{22}^{N-Q-1}\}$. $\therefore \forall p \geq N - Q - 1$, $\text{rank}(\mathcal{O}_{p+2}) = \text{rank}(\mathcal{O}_{p+1})$. Correspondingly the minimal number of previous steps that satisfies the rank test can be defined as the minimal integer that satisfies the $\text{rank}(\mathcal{O}_{p+2}) = \text{rank}(\mathcal{O}_{p+1})$.

$$(5.16) \quad p_* = \min \Pi_{\mathcal{O}} - 1,$$

where

$$(5.17) \quad \Pi_{\mathcal{O}} = \{l \in \mathbb{N}^+, r_{\mathcal{O}}(l) = r_{\mathcal{O}}(l + 1)\}.$$

Because of the monotonicity of integer function $r_{\mathcal{O}}(\cdot)$, p_* can be found in a sequential sense. \blacksquare

The fact that one can determine the closure dynamics of any nonlinear system with dual linear closure given *all* previous resolved states is not surprising. It will be shown shortly that this is possible for a slightly more general case.

Proposition 5.7. *Closure dynamics of any nonlinear dynamical system with dual linear closure can be determined as a linear combination of $\hat{\mathbf{x}}^{n-1}, \dots, \hat{\mathbf{x}}^{n-p}, \boldsymbol{\delta}^n, \dots, \boldsymbol{\delta}^{n-p}$, with $p \leq N - Q - 1$.*

As a trivial observation, if we replace closure with $\tilde{\mathbf{x}}^n$, one can easily obtain the following rank test as a criterion.

Proposition 5.8. *For a nonlinear dynamical system with dual linear closure, if*

$$\text{rank}(V_p) = \text{rank}\left(\begin{bmatrix} C^\top \\ V_p^\top \end{bmatrix}\right),$$

where $C^\top = [I_{N-Q \times N-Q} \quad \mathbf{0}]$, then $\tilde{\mathbf{x}}^n$ is observable from a p time delayed observable space.

The key ingredient of [Theorem 5.2](#) is the exploitation of *projection* equations in the *dual linear* closure setting, which may be overlooked since they share the same information as previous observables in the statistical sense if the initial condition is fully known. Without the explicit usage of *projection* equations, one can obtain a closure with explicit memory dependence on *all* previous observables, but is correspondingly applicable to a more general system stated in [Definition 5.9](#).

Definition 5.9 (Nonlinear dynamical system with linear closure). *A nonlinear dynamical system with linear closure is defined as*

$$(5.18) \quad \frac{d}{dt} \begin{bmatrix} \hat{\mathbf{x}} \\ \tilde{\mathbf{x}} \end{bmatrix} = \begin{bmatrix} \mathbf{F}(\hat{\mathbf{x}}, \tilde{\mathbf{x}}) \\ \mathbf{H}(\hat{\mathbf{x}}) + A_{22}\tilde{\mathbf{x}} \end{bmatrix},$$

where $\hat{\mathbf{x}} \in \mathbb{R}^Q$, $\tilde{\mathbf{x}} \in \mathbb{R}^{N-Q}$, $\mathbf{F}(\cdot) : \mathbb{R}^N \mapsto \mathbb{R}^Q$ and $\mathbf{H}(\cdot) : \mathbb{R}^Q \mapsto \mathbb{R}^{N-Q}$, $A_{22} \in \mathbb{R}^{N-Q \times N-Q}$ with $\boldsymbol{\delta} = A_{12}\tilde{\mathbf{x}}$ and $Q \in \mathbb{N}$, $Q < N$.

Corollary 5.10. *With only evolution equations, one can write the following equation for a first order forward discretized dynamical system $\forall n, p \in \mathbb{N}$, $n > p$*

$$\tilde{\mathbf{x}}^n = (I + \Delta t A_{22})^p \tilde{\mathbf{x}}^{n-p} + \sum_{l=0}^{p-1} \Delta t (I + \Delta t A_{22})^l \mathbf{H}(\hat{\mathbf{x}}^{n-l-1}),$$

which links the unresolved states between any two time instances.

Proof. Considering only the *evolution* equations, one can write the following in matrix form

$$(5.19) \quad \mathbf{\Gamma}_p^e \mathbf{X}_p = \boldsymbol{\Sigma}_p^e,$$

where

$$(5.20) \quad \mathbf{\Gamma}_p^e = \begin{bmatrix} I & -(I + \Delta t A_{22}) & \dots & \mathbf{0} \\ \mathbf{0} & I & -(I + \Delta t A_{22}) & \dots \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \dots & I & -(I + \Delta t A_{22}) \end{bmatrix},$$

$$(5.21) \quad \mathbf{X}_p = \begin{bmatrix} \tilde{\mathbf{x}}^n \\ \tilde{\mathbf{x}}^{n-1} \\ \vdots \\ \vdots \\ \vdots \\ \tilde{\mathbf{x}}^{n-p} \end{bmatrix}, \quad \Sigma_p^e = \begin{bmatrix} \Delta t \mathbf{H}(\hat{\mathbf{x}}^{n-1}) \\ \Delta t \mathbf{H}(\hat{\mathbf{x}}^{n-2}) \\ \vdots \\ \Delta t \mathbf{H}(\hat{\mathbf{x}}^{n-p}) \end{bmatrix}.$$

Recall we are interested in $\hat{\mathbf{x}}^n$, with several row operations on the first row block,

$$(5.22) \quad \Gamma_p^e \rightarrow \begin{bmatrix} I & \mathbf{0} & \dots & -(I + \Delta t A_{22})^p \\ \dots & \dots & \dots & \dots \end{bmatrix},$$

and correspondingly

$$(5.23) \quad \Sigma_p^e \rightarrow \begin{bmatrix} \sum_{l=0}^{p-1} \Delta t (I + \Delta t A_{22})^l \mathbf{H}(\hat{\mathbf{x}}^{n-l-1}) \\ \vdots \end{bmatrix},$$

Therefore,

$$(5.24) \quad \tilde{\mathbf{x}}^n = (I + \Delta t A_{22})^p \tilde{\mathbf{x}}^{n-p} + \sum_{l=0}^{p-1} \Delta t (I + \Delta t A_{22})^l \mathbf{H}(\hat{\mathbf{x}}^{n-l-1}). \quad \blacksquare$$

The implication is that, if $\tilde{\mathbf{x}}$ is known at one previous time instant, the future of $\tilde{\mathbf{x}}$ starting from that point is completely determined by $\hat{\mathbf{x}}$ in a convolutional sense. For example, starting from the initial condition, we have the following result often seen in linear systems theory:

Proposition 5.11. *For a nonlinear dynamical system with linear closure, if $\tilde{\mathbf{x}}^0$ is known, one can uniquely determine $\tilde{\mathbf{x}}^n$ in the following*

$$(5.25) \quad \tilde{\mathbf{x}}^n = (I + \Delta t A_{22})^n \tilde{\mathbf{x}}^0 + \sum_{l=0}^{n-1} \Delta t (I + \Delta t A_{22})^l \mathbf{H}(\hat{\mathbf{x}}^{n-l-1}).$$

The present framework exploits the fact that although the closure is explicitly based on *all* previous information of the observables, the operator driving this function might only possess a finite memory dependence as indicated in [Theorem 5.2](#). Therefore, the essential structure of the closure may be much more compact.

6. Results - Non-linear systems. In this section, the operator inference framework is used to derive closures for several problems ranging from chaotic and non-chaotic nonlinear ordinary differential equation (ODE) systems to nonlinear partial differential equations (PDE).

6.1. Van del Pol system.

6.1.1. Problem description. The Van del Pol (VdP) system with first order forward discretization is

$$(6.1) \quad \begin{bmatrix} x_1^{n+1} \\ x_2^{n+1} \end{bmatrix} = \begin{bmatrix} x_1^n \\ x_2^n \end{bmatrix} + \Delta t \begin{bmatrix} x_2^n \\ \mu(1 - x_1^n x_1^n)x_2^n - x_1^n \end{bmatrix},$$

where $\mu = 2$, $\hat{x}(0) = x_1^0 = 1$, $\tilde{x}(0) = x_2^0 = 0$, and $\delta = x_2$. The simulation is run from $t \in [0, 60]$ collecting $\{\hat{x}, \delta\}$ as data over 6000 snapshots with a $\Delta t = 0.01$. The first 30% of data is set as training data and the rest is set for testing.

Consider $\hat{x} = x_1$, $\tilde{x} = x_2$ thus $N = 2$ and $Q = 1$. Correspondingly, the ROM formulation is given below with linear superposition of multi-time effects assumption

$$(6.2) \quad \hat{x}^{n+1} = \hat{x}^n + \Delta t \delta^n,$$

$$(6.3) \quad \delta^{n+1} = \delta^n + \Delta t \sum_{i=0}^p G_i.$$

For VdP system, the exact solution for the closure dynamics with $p = 0$ is

$$(6.4) \quad \frac{\delta^{n+1} - \delta^n}{\Delta t} = \mu(1 - \hat{x}^n \hat{x}^n) \delta^n - \hat{x}^n = -\hat{x}^n + 2\delta^n - 2\hat{x}^n \hat{x}^n \delta^n.$$

6.1.2. Model selection for polynomial regression. To determine the underlying sparse dynamics, the *lasso* path is computed and presented in Figures 8 and 9. It can be seen that an elbow is present in the error plot as λ near 10^{-10} , where a number of non-zero terms jump above 3 to 9, causing a slight increase in MSE. Thus, the optimal λ is chosen as 10^{-10} according to the Pareto front.

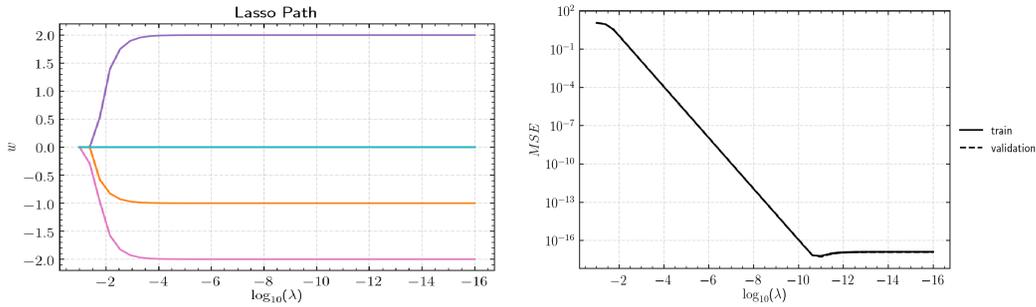


Figure 8. lasso path for 2D VdP system. Left: coefficients. Right: MSE.

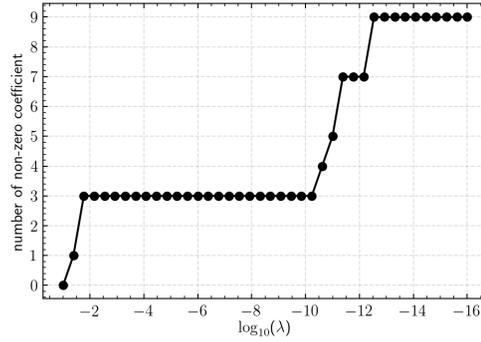


Figure 9. lasso path for 2D VdP system: number of non-zero terms

6.1.3. A posteriori evaluation of model performance. The corresponding model performance in an a posteriori sense for both training and testing data is excellent, as shown below in Figure 10.

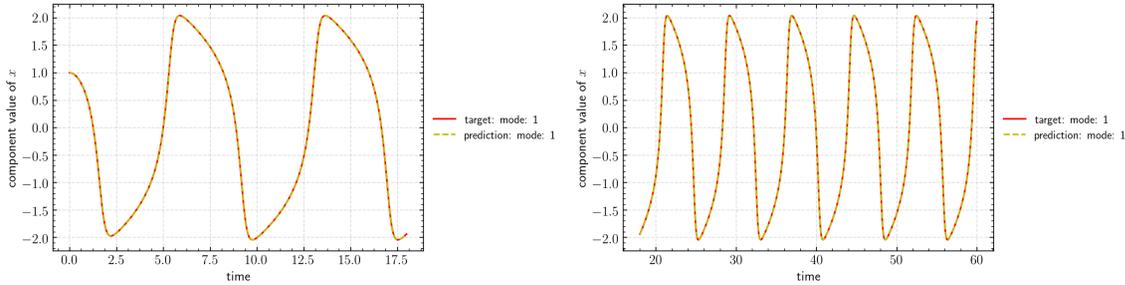


Figure 10. A posteriori model performance on 2D VdP system. Left: training data. Right: testing data.

6.2. Duffing Map.

6.2.1. Problem description. The Duffing map is a classic example of a chaotic map. We take the form

$$(6.5) \quad x_1^{n+1} = x_1^n + \Delta t(x_2^n - x_1^n),$$

$$(6.6) \quad x_2^{n+1} = x_2^n + \Delta t(-bx_1^n + (a-1)x_2^n - (x_2^n)^3),$$

with $a = 2.75$ and $b = 0.2$, $\Delta t = 1$, $x_1(0) = x_1^0 = 0.5$, $x_2(0) = x_2^0 = 0$. The resolved variable $\hat{x} = x_1$. We simulate this system up to 6000 steps with the first 30% for training, and the rest for testing. For this case, the corresponding closure dynamics for δ is

$$(6.7) \quad \delta^{n+1} = a\delta^n - (\delta^n)^3 - b\hat{x}^n.$$

6.2.2. Model selection. As displayed in Figures 11 and 12, by sweeping λ , an optimal value of $\lambda = 10^{-10}$ is found. At that sparsity level, the resulting expression is given as follows:

$$(6.8) \quad \delta^{n+1} = \delta^n + \Delta t(-0.199999\hat{x}^n + 1.749999\delta^n - 0.999999\delta^{n3} + 2.99 \times 10^{-11}\hat{x}^{n2} - 1.34 \times 10^{-8}\hat{x}^{n3} + 2.41 \times 10^{-11}\delta^n\hat{x}^n + 6.81 \times 10^{-10}\delta^n\hat{x}^{n2}).$$

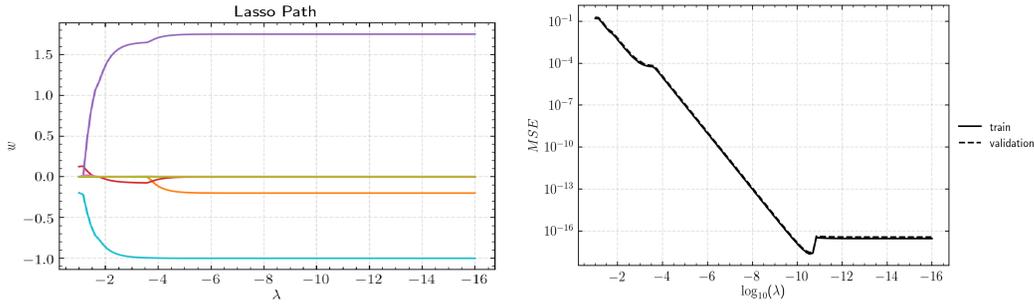


Figure 11. lasso path for 2D Duffing system. Left: coefficient. Right: MSE.

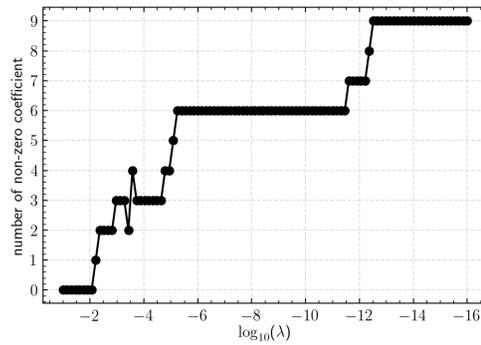


Figure 12. lasso path for 2D Duffing system: number of non-zero terms

6.2.3. A posteriori evaluation of model performance. For a chaotic system, since it is extremely difficult to achieve accurate long time predictions, models are most often evaluated in a variety of ways. These include subjective visual inspection [21] or measures for the attractor [2] such as maximum Lyapunov exponent [17], correlation dimension and other time averaged characteristics [28]. The first approach, although perhaps the most widely used [30][39][49], can sometimes be misleading [15].

In this work, we first show there is excellent correspondence in maximum Lyapunov exponent (MLE) and correlation dimension γ , computed on both the ground truth time series and modeled time series for both training and testing data. Following this, we employ a null hypothesis test proposed by Diks et al. [15] to show that the attractor reconstructed by embedding the time series predicted by our model is indeed close to the phase space reconstruction of the ground truth within a confidence interval. As suggested by Diks, the null hypothesis that the two delay vectors are drawn from the same multidimensional probability distribution is accepted if $S < 3$. Diks criterion has been previously employed as an early stop criterion during the training of neural networks [3].

The comparison of the predicted time series between the modeled system and ground truth is displayed in Figure 13 for training and testing data. Due to the chaotic nature of the dynamics, direct measurement of the MSE is not suitable for this case. Examination of the MLE and correlation dimension in Table 1 shows excellent agreement. Furthermore,

Diks test shows $|S| = 1.003$ for training data and $|S| = 1.588$ for testing data, which further confirms the validity of the model. Details of the implementation of Diks criterion are given in [Appendix B](#).

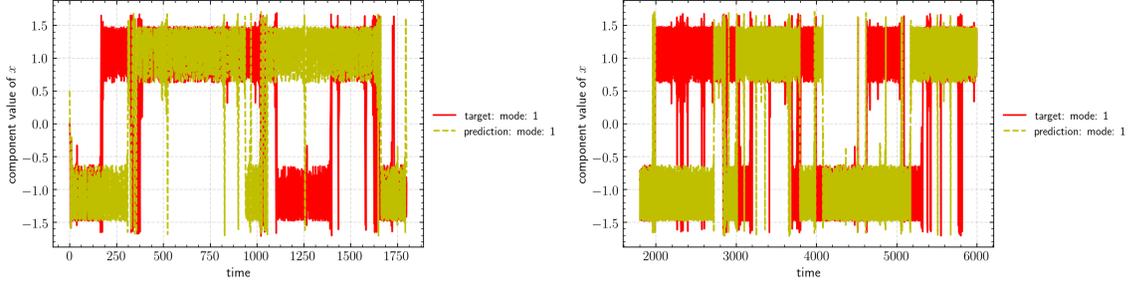


Figure 13. *A posteriori model performance on the Duffing map. Left: training data. Right: testing data.*

Table 1

Comparison of MLE and correlation dimension between truth and model

case	MLE	γ
true train	0.98	1.12
model train	0.97	1.12
true test	0.97	1.13
model test	0.97	1.13

6.3. Lorenz system.

6.3.1. Problem description. The corresponding first order forward discretized Lorenz system is given as:

$$(6.9) \quad \begin{bmatrix} x_1^{n+1} \\ x_2^{n+1} \\ x_3^{n+1} \end{bmatrix} = \begin{bmatrix} x_1^n + \Delta t \sigma (x_2^n - x_1^n) \\ x_2^n + \Delta t (x_1^n (\rho - x_3^n) - x_2^n) \\ x_3^n + \Delta t (x_1^n x_2^n - \beta x_3^n), \end{bmatrix}$$

with $x_1(0) = 0.5$, $x_2(0) = x_3(0) = 0$. Parameters for each case are shown in [Table 2](#) where the only difference is ρ .

Table 2

Parameters of Lorenz system for chaotic and nonchaotic cases

case	σ	β	ρ
non-chaotic	10	8/3	15
chaotic	10	8/3	35

For the non-chaotic case, the simulation time is $t = [0, 20]$ with 8000 snapshots; and for the chaotic case, the simulation time is $t = [0, 400]$ with 40000 snapshots. The snapshots are equally split between training and testing sets.

For the Lorenz system with $\hat{x} = x_1$, $\delta = \sigma x_2$, one can find the analytical closure for δ with $p = 1$ after some algebra:

$$(6.10) \quad \delta^{n+1} = (1 - \Delta t)\delta^n + \sigma\Delta t\hat{x}^n \left((1 - \beta\Delta t) \left(\frac{\delta^n + (\Delta t - 1)\delta^{n-1}}{\sigma\hat{x}^{n-1}\Delta t} \right) - \frac{\hat{x}^{n-1}\delta^{n-1}}{\sigma} + \rho\beta\Delta t \right),$$

which clearly involves *cross time* features and *rational* forms instead of pure polynomial forms.

The corresponding ROM formulation is given as

$$(6.11) \quad \hat{x}^{n+1} = \hat{x}^n - \Delta t\sigma\hat{x}^n + \Delta t\delta^n,$$

$$(6.12) \quad \delta^{n+1} = \delta^n + \Delta t\mathbf{G}(\hat{x}^n, \hat{x}^{n-1}, \delta^n),$$

where \mathbf{G} is modeled by a neural network.

Standard polynomial regression is found to be unsuitable to extract governing equations in this case. A recently developed method called implicit-SINDy [29], which can account for non-rational functions could perhaps improve predictions. Alternatively, we employ an artificial neural network model with $p = 1$ and consider full memory interaction between different time instances. The architecture of the neural network is chosen as 4-16-16-1 for both chaotic and non-chaotic cases with $p = 1$ and *tanh* as the activation function. The neural network model is trained for 16000 epochs with the Adam optimizer with a learning rate of 0.0001, a mini-batch size of 256 and the last 10% of training data is split as validation set to monitor generalization.

6.3.2. A posteriori evaluation of model performance. For the non-chaotic case, the model performs well for both training and testing data, as shown in Figure 14.

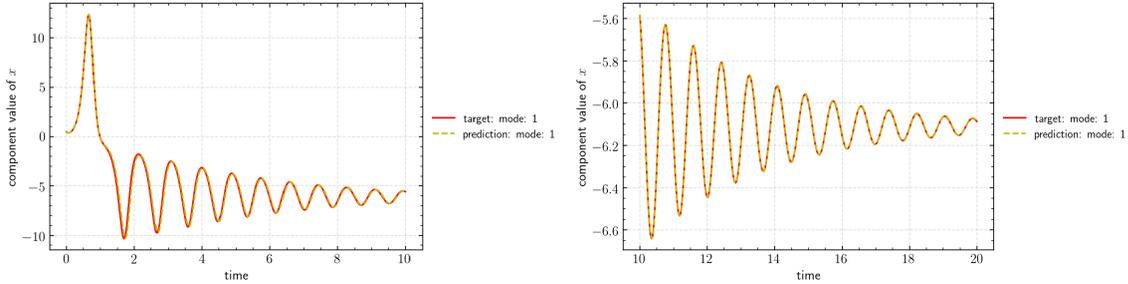


Figure 14. A posteriori model performance on the non-chaotic Lorenz system. Left: training data. Right: testing data.

For the chaotic case, results are shown in Figure 15 for training and testing evaluations. Table 3 shows that both MLE and correlation dimension are in accordance with the truth. Furthermore, Diks criterion shows $|S| = 1.509$ for training data and $|S| = 2.83$ for testing data, which implies that the null hypothesis is accepted. Details of implementation are provided in Appendix B.

Table 3

Comparison of MLE and correlation dimension between truth and model

case	MLE	γ
true train	0.044	1.34
model train	0.042	1.33
true test	0.041	1.34
model test	0.041	1.34

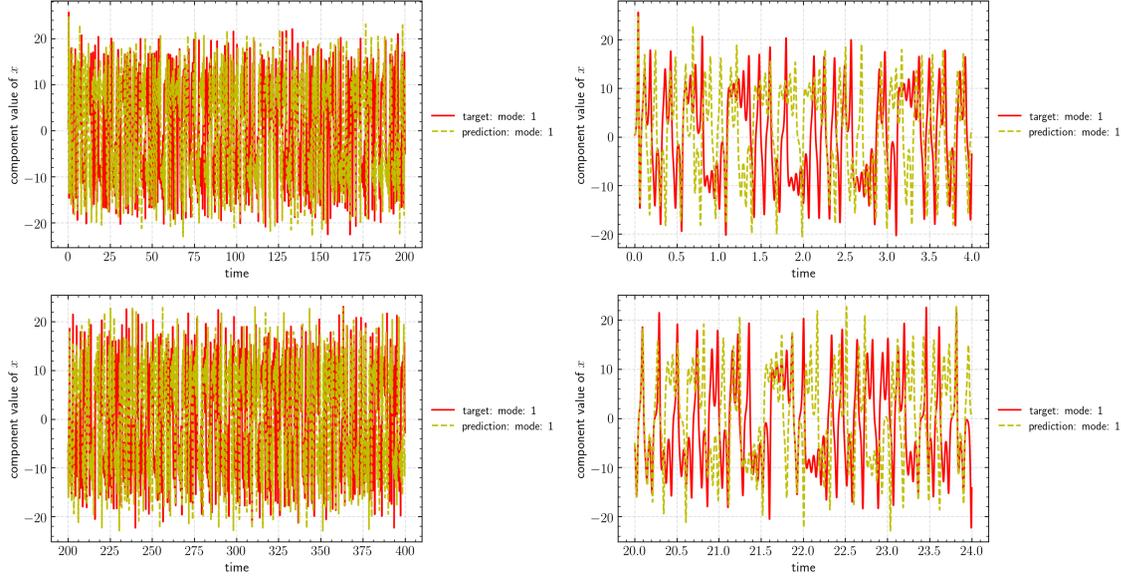


Figure 15. A posteriori model performance on the chaotic Lorenz system. Top left: training data. Bottom left: testing data. Top right: training data zoomed in $t \in [0, 4]$. Bottom right: testing data zoomed in $t \in [20, 24]$

6.4. One dimensional viscous Burgers equation.

6.4.1. Problem description. In this section, the one dimensional viscous Burgers equation is considered in a periodic domain $x \in [0, 2\pi]$

$$(6.13) \quad \frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} = \nu \frac{\partial^2 u}{\partial x^2}$$

with $\nu = 0.02$, $t \in [0, 10]$, $u(x, 0) = \sin(x)$.

Using a standard pseudo-spectral method with two-thirds dealiasing and with Runge-Kutta 3rd order SSP scheme for time stepping, the system is resolved with 1024 grid points uniformly distributed in space, and a time step $\Delta t = 0.01\Delta x$. 2000 snapshots of $u(x, t)$ are uniformly collected in time. For the setup of coarse graining, we use a spectral filter to obtain the state and corresponding closure with wavenumber k ranging from -6 to 5 . The

corresponding equation in spectral form for k^{th} wavenumber or mode is

$$(6.14) \quad \frac{d\hat{u}_k}{dt} = -\nu k^2 \hat{u}_k - \frac{ik}{2} \sum_{p+q=k} \hat{u}_p \hat{u}_q = -\nu k^2 \hat{u}_k - \frac{ik}{2} \sum_{p+q=k, p \in F, q \in F} \hat{u}_p \hat{u}_q \\ - \frac{ik}{2} \sum_{p+q=k, p \in F, q \in G} \hat{u}_p \hat{u}_q - \frac{ik}{2} \sum_{p+q=k, p \in G, q \in F} \hat{u}_p \hat{u}_q - \frac{ik}{2} \sum_{p+q=k, p \in G, q \in G} \hat{u}_p \hat{u}_q,$$

where $u(x, 0) = \sin(x)$; $x \in [0, 2\pi]$; $k \in \{-N/2, \dots, N/2 - 1\}$; index set of resolved modes $F = \{-Q/2, \dots, Q/2 - 1\}$; index set of unresolved modes: $G = \{-N/2, \dots, -Q/2 - 1, Q/2, \dots, N/2 - 1\}$. The closure is the sum of last three terms in (6.14), and noticing there is a symmetry in the solution with sine wave initial condition, a truncation corresponding to $Q = 6$ is considered. Only the imaginary part of \hat{u}_k with k ranges from -6 to -1 is considered. For the evaluation of the closure model, we consider $0 \leq t \leq 4$ as our training data and $4 < t \leq 10$ as testing data.

6.4.2. Model selection. For the polynomial model, the optimal time delay p and polynomial order k is chosen by sweeping p from 0 to 2. For each p , the optimal k and corresponding λ is extracted.

For the application of the ANN, the best model is selected from a range of hyperparameters with p ranging from 0 to 2. Two hidden layers are used with identical numbers of hidden units for each layer as 4, 8, 12, 16 and type of activation as ReLU, SeLU, tanh. The optimal model was chosen as that which yields the most satisfactory validation result with the smallest number of parameters. We found this to be $p = 2$, with 12 hidden units and the tanh activation function. The type of activation function does not appear to be critical in this case, which may be a consequence of the fact that it is not a deep neural network where the vanishing gradient problem may be significant [19].

6.4.3. A posteriori evaluation of model performance. As seen in Figures 16 and 18, both SINDy and ANN perform well on the training data. When evaluated against unseen testing data, however, performance of SINDy was seen to deteriorate as displayed in Figure 17 as a consequence of the extrapolation going out of bounds with high order polynomial features. The ANN is a convergent series of infinite polynomials, and therefore the corresponding model remains less unbounded compared to the polynomial model as shown in Figure 18.

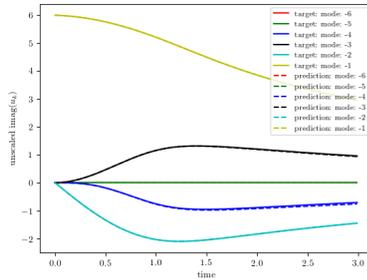


Figure 16. A posteriori model performance on training data : 1D VBE using polynomial closures.

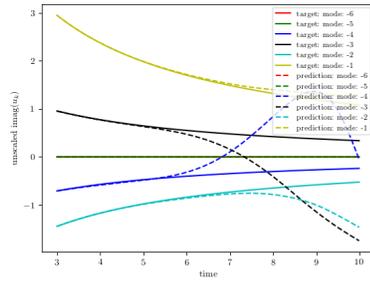


Figure 17. A posteriori model performance on testing data : 1D VBE using polynomial closures.

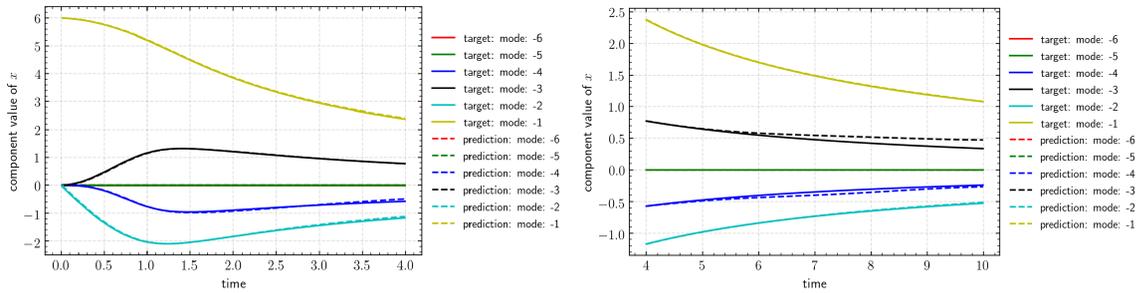


Figure 18. A posteriori model performance on 1D VBE with ANN. Left: training data. Right: testing data.

Comparison of the results on unseen testing data with ANN at snapshots $t = 4.5$, $t = 6.0$, $t = 7.5$ and $t = 9.0$ between the data-driven model, no closure, and ground truth in physical space is shown in Figure 19. The results highlight the importance of the closure in predicting the future state of this system. The ANN model performs particularly well between $t \in [4, 6]$, with a slight degradation in performance at later times.

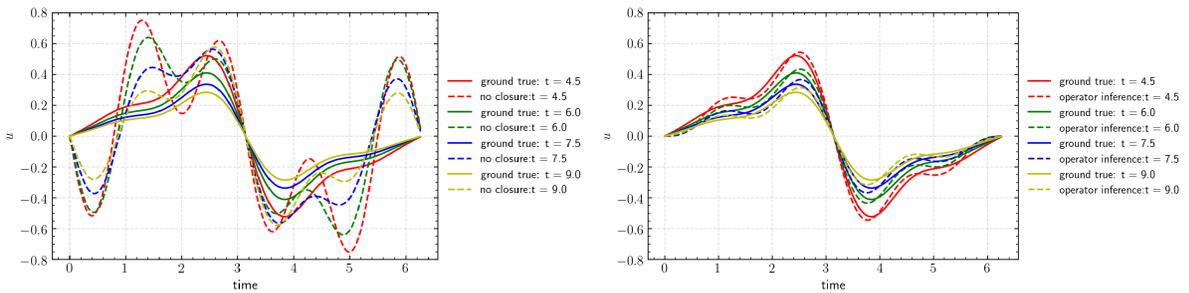


Figure 19. A posteriori model performance on 1D VBE with ANN. Left: without closure. Right: with operator inference closure.

7. Conclusion. An operator inference framework was presented, with the goal of developing closures for reduced models of dynamical systems². Dynamic memory is embedded into the equations and the evolution of this term is parametrized via polynomial features and artificial neural networks (ANN). The polynomial model is determined using non-linear regression and *lasso* with Pareto-front-based model selection. The ANN model is determined using gradient-based methods with weight decay regularization. By assuming that different time instances are decoupled from each other, the exponential growth of the number of parameters is limited to a linear growth. For special types of non-linear systems, the closure dynamics was proven to have a compact memory, and the form of the closure is shown to be precisely discoverable using a sparse set of features. Numerical evaluations of the model on non-chaotic and chaotic dynamical systems are used to evaluate the viability of the procedure, with an emphasis on model selection and a posteriori prediction of unseen data.

Acknowledgements. This work was supported by AFOSR and AFRL under grants FA9550-16-1-0309 & FA9550-17-1-0195. The authors thank Mr. Sven Giorno for numerical experimentation, Prof. Cees Diks for discussion on implementing the Diks criterion, and Ms. Helen Zhang for comments on the manuscript.

Appendix A. Comparison with Elman's recurrent neural network.

Elman's network[16] is one of the earliest [48] recurrent neural network models, and was originally proposed to represent temporal structure in linguistics. Although Elman's network is similar to a standard feedforward neural network (FNN), the key difference is that its input includes an additional feedback, and thus the memory effect is addressed in a lossy sense [19] using one previous step.

In this section, we will highlight similarities and differences between the operator inference framework for closure modeling (3.1) and (3.2) and Elman's model. Given a general predictive task for a discrete dynamical system: $\{x_i\}_{i=1,\dots}$, $x_i \in \mathbb{R}^Q$, $i \in \mathbb{N}^+$, $Q \in \mathbb{N}^+$, Elman's network is:

$$(A.1) \quad x_{i+1} = \mathcal{C}(h_{i+1}),$$

$$(A.2) \quad h_{i+1} = \mathcal{H}(h_i, x_i),$$

where $\mathcal{H}(\cdot)$ and $\mathcal{C}(\cdot)$ are perceptrons and $h_i \in \mathbb{R}^H$, $H \in \mathbb{N}^+$ is the number of units in the context layer. On the other hand, if one considers a simplified discrete case of (3.1) and (3.2) with $p = 0$, following the same notation, one has:

$$(A.3) \quad x_{i+1} = f(x_i) + h_i,$$

$$(A.4) \quad h_{i+1} = G(h_i, x_i),$$

where $f : \mathbb{R}^Q \mapsto \mathbb{R}^Q$ is known, while $G : \mathbb{R}^Q \times \mathbb{R}^Q \mapsto \mathbb{R}^Q$ is unknown. The similarity is that both (A.2) and (A.4) address the memory effect and extract the dynamics in the same fashion. However, there are at least three different aspects:

²Sample code available at: https://github.com/pswpswpsw/siads_data_driven_closure.git

- Elman’s network assumes output dependence only on newly activated hidden units h_{i+1} , while our model at $p = 0$ considers output dependence on previously activated hidden units h_I , together with the current input x_I . Our model also extends the case to $p > 0$,
- Our model decouples the evolution processes of hidden units and states while Elman’s is formulated in a sequential fashion,
- Elman’s network requires the determination of all relationships, i.e., the perceptrons, in a purely data-driven fashion, whereas the structure of state evolution is considered known in our operator inference framework.

Appendix B. Implementation of Diks criterion.

Diks et al. [15] developed a test that evaluates whether two attractors are similar. Diks criterion follows statistical inference and can provide probabilistic confidence bounds. In our work, this criterion is used to compare the reconstructed dynamics of an attractor with the ground truth. The method is based on testing a null hypothesis: *two sets of delay vectors are drawn from the same multidimensional probability distribution*. It was later employed by Bakker as a monitoring metric *during* the training of ANNs for time series modeling. The time series is divided into segments of length l and averaged. To cope with the fractal probability distribution of the chaotic attractor, smoothing is performed via a Gaussian kernel. A bandwidth d is determined by performing sweeps on another trajectory and choosing the d that reveals the highest discrepancy between the two time series. Other hyperparameters are the embedding dimension m , and delay time τ . τ is chosen as the first local minimum of mutual information of Fraser, and m is simply chosen as 2 for the Duffing map and 3 for the Lorenz system.

Given two sets $\{\mathbf{X}_i\}_{i=1}^{N_1}$ and $\{\mathbf{Y}_i\}_{i=1}^{N_2}$ and realizations $\{\mathbf{x}_i\}_{i=1}^{N_1}$ and $\{\mathbf{y}_i\}_{i=1}^{N_2}$, the square root of Q defines a distance between the two probability distribution of delay vectors. \hat{Q} is an unbiased estimator of Q and given by:

$$(B.1) \quad \hat{Q} = \frac{1}{\binom{N_1}{2}} \sum_{1 \leq i < j \leq N_1} h(\mathbf{X}_i, \mathbf{X}_j) + \frac{1}{\binom{N_2}{2}} \sum_{1 \leq i < j \leq N_2} h(\mathbf{Y}_i, \mathbf{Y}_j) - \frac{2}{N_1 N_2} \sum_{i=1}^{N_1} \sum_{j=1}^{N_2} h(\mathbf{X}_i, \mathbf{Y}_j).$$

The variance of \hat{Q} under a null hypothesis and conditionally on the set of $N = N_1 + N_2$ observed vectors is given by:

$$(B.2) \quad V_c(\hat{Q}) = \frac{2(N-1)^2(N-2)}{N_1(N_1-1)N_2(N_2-1)(N-3)} \frac{1}{\binom{N}{2}} \sum_{1 \leq i < j \leq N} \phi_{ij}^2,$$

where

$$\begin{aligned} \phi_{ij} &= H_{ij} - g_i - g_j, \\ h(\mathbf{s}, \mathbf{t}) &= e^{-|\mathbf{s}-\mathbf{t}|/4d^2}, \end{aligned}$$

and

$$H_{ij} = h(\mathbf{z}_i, \mathbf{z}_j) - \frac{1}{\binom{N}{2}} \sum_{1 \leq i < j \leq N} h(\mathbf{z}_i, \mathbf{z}_j),$$

and $g_i = \frac{1}{N-2} \sum_{j, j \neq i} H_{ij}$, where z_i is defined as

$$z_i = \begin{cases} \mathbf{x}_i, & \text{for } 1 \leq i \leq N_1 \\ \mathbf{y}_{i-N_1}, & \text{for } N_1 < i \leq N \end{cases}.$$

Note that $S = \hat{Q}/V_c(\hat{Q})$ is a random variable with zero mean and unit standard derivation under the null hypothesis. As suggested by Diks, we reject the null hypothesis with more than 95% confidence for $S > 3$.

In this work, for the Duffing map, the optimal $d = 0.0001$, l is chosen as 100, τ is chosen as 20. For the Lorenz system, the optimal $d = 0.001$ and l is chosen as 100, τ is chosen as 25.

REFERENCES

- [1] M. ABADI, A. AGARWAL, P. BARHAM, E. BREVDO, Z. CHEN, C. CITRO, G. S. CORRADO, A. DAVIS, J. DEAN, M. DEVIN, S. GHEMAWAT, I. J. GOODFELLOW, A. HARP, G. IRVING, M. ISARD, Y. JIA, R. JÓZEFOWICZ, L. KAISER, M. KUDLUR, J. LEVENBERG, D. MANÉ, R. MONGA, S. MOORE, D. G. MURRAY, C. OLAH, M. SCHUSTER, J. SHLENS, B. STEINER, I. SUTSKEVER, K. TALWAR, P. A. TUCKER, V. VANHOUCHE, V. VASUDEVAN, F. B. VIÉGAS, O. VINYALS, P. WARDEN, M. WATTENBERG, M. WICKE, Y. YU, AND X. ZHENG, *Tensorflow: Large-scale machine learning on heterogeneous distributed systems*, CoRR, abs/1603.04467 (2016), <http://arxiv.org/abs/1603.04467>.
- [2] L. A. AGUIRRE AND C. LETELLIER, *Modeling nonlinear dynamics and chaos: a review*, Mathematical Problems in Engineering, 2009 (2009).
- [3] R. BAKKER, J. C. SCHOUTEN, C. L. GILES, F. C. TAKENS, AND C. M. VAN DEN BLEEK, *Learning chaotic attractors by neural networks*, Neural Comput., 12 (2000), pp. 2355–2383, <https://doi.org/10.1162/089976600300014971>, <http://dx.doi.org/10.1162/089976600300014971>.
- [4] C. A. BEATTIE AND S. GUGERCIN, *Krylov-based model reduction of second-order systems with proportional damping*, in Decision and Control, 2005 and 2005 European Control Conference. CDC-ECC'05. 44th IEEE Conference on, IEEE, 2005, pp. 2278–2283.
- [5] M. BENOSMAN, J. BORGGAARD, O. SAN, AND B. KRAMER, *Learning-based robust stabilization for reduced-order models of 2d and 3d boussinesq equations*, Applied Mathematical Modelling, 49 (2017), pp. 162–181.
- [6] G. BERKOOZ, P. HOLMES, AND J. L. LUMLEY, *The proper orthogonal decomposition in the analysis of turbulent flows*, Annual review of fluid mechanics, 25 (1993), pp. 539–575.
- [7] S. A. BILLINGS, *Nonlinear system identification: NARMAX methods in the time, frequency, and spatio-temporal domains*, John Wiley & Sons, 2013.
- [8] S. L. BRUNTON, J. L. PROCTOR, AND J. N. KUTZ, *Discovering governing equations from data by sparse identification of nonlinear dynamical systems*, Proceedings of the National Academy of Sciences, 113 (2016), pp. 3932–3937.
- [9] K. CARLBERG, C. FARHAT, J. CORTIAL, AND D. AMSALLEM, *The gnat method for nonlinear model reduction: effective implementation and application to computational fluid dynamics and turbulent flows*, Journal of Computational Physics, 242 (2013), pp. 623–647.
- [10] F. CHINESTA, P. LADEVEZE, AND E. CUETO, *A short review on model order reduction based on proper generalized decomposition*, Archives of Computational Methods in Engineering, 18 (2011), p. 395.
- [11] F. CHOLLET ET AL., *Keras*. <https://keras.io>, 2015.
- [12] A. J. CHORIN AND O. H. HALD, *Estimating the uncertainty in underresolved nonlinear dynamics*, Mathematics and Mechanics of Solids, 19 (2014), pp. 28–38.
- [13] A. J. CHORIN, O. H. HALD, AND R. KUPFERMAN, *Optimal prediction with memory*, Physica D: Nonlinear Phenomena, 166 (2002), pp. 239–257.
- [14] M. COUPLET, C. BASDEVANT, AND P. SAGAUT, *Calibrated reduced-order pod-galerkin system for fluid flow modelling*, Journal of Computational Physics, 207 (2005), pp. 192–220.

- [15] C. DIKS, W. VAN ZWET, F. TAKENS, AND J. DEGOEDE, *Detecting differences between delay vector distributions*, Physical Review E, 53 (1996), p. 2169.
- [16] J. L. ELMAN, *Finding structure in time*, Cognitive science, 14 (1990), pp. 179–211.
- [17] R. GENÇAY AND T. LIU, *Nonlinear modelling and prediction with feedforward and recurrent networks*, Physica D: Nonlinear Phenomena, 108 (1997), pp. 119–134.
- [18] M. GERMANO, U. PIOMELLI, P. MOIN, AND W. H. CABOT, *A dynamic subgrid-scale eddy viscosity model*, Physics of Fluids A: Fluid Dynamics, 3 (1991), pp. 1760–1765.
- [19] I. GOODFELLOW, Y. BENGIO, AND A. COURVILLE, *Deep learning*, MIT press, 2016.
- [20] A. GOUASMI, E. J. PARISH, AND K. DURAISAMY, *A priori estimation of memory effects in reduced-order models of nonlinear systems using the mori–zwanzig formalism*, in Proc. R. Soc. A, vol. 473, The Royal Society, 2017, p. 20170385.
- [21] M. HAN, J. XI, S. XU, AND F.-L. YIN, *Prediction of chaotic time series based on the recurrent predictor neural network*, IEEE transactions on signal processing, 52 (2004), pp. 3409–3416.
- [22] P. HOLMES, *Turbulence, coherent structures, dynamical systems and symmetry*, Cambridge university press, 2012.
- [23] K. HORNIK, M. STINCHCOMBE, AND H. WHITE, *Multilayer feedforward networks are universal approximators*, Neural Networks, 2 (1989), pp. 359–366, [https://doi.org/10.1016/0893-6080\(89\)90020-8](https://doi.org/10.1016/0893-6080(89)90020-8), <https://arxiv.org/abs/arXiv:1011.1669v3>.
- [24] T. J. HUGHES, G. R. FEIJÓO, L. MAZZEI, AND J.-B. QUINCY, *The variational multiscale method a paradigm for computational mechanics*, Computer methods in applied mechanics and engineering, 166 (1998), pp. 3–24.
- [25] R. IBÁÑEZ, E. ABISSET-CHAVANNE, J. V. AGUADO, D. GONZALEZ, E. CUETO, AND F. CHINESTA, *A manifold learning approach to data-driven computational elasticity and inelasticity*, Archives of Computational Methods in Engineering, 25 (2018), pp. 47–57.
- [26] D. KINGMA AND J. BA, *Adam: A method for stochastic optimization*, arXiv preprint arXiv:1412.6980, (2014).
- [27] V. KOUZNETSOVA, W. BREKELMANS, AND F. BAAIJENS, *An approach to micro-macro modeling of heterogeneous materials*, Computational Mechanics, 27 (2001), pp. 37–48.
- [28] H. LIN, W. CHEN, AND A. TSUTSUMI, *Long-term prediction of nonlinear hydrodynamics in bubble columns by using artificial neural networks*, Chemical Engineering and Processing: Process Intensification, 42 (2003), pp. 611–620.
- [29] N. M. MANGAN, S. L. BRUNTON, J. L. PROCTOR, AND J. N. KUTZ, *Inferring biological networks by sparse identification of nonlinear dynamics*, IEEE Transactions on Molecular, Biological and Multi-Scale Communications, 2 (2016), pp. 52–63.
- [30] T. MIYOSHI, H. ICHIHASHI, S. OKAMOTO, AND T. HAYAKAWA, *Learning chaotic dynamics in recurrent rbf network*, in Neural Networks, 1995. Proceedings., IEEE International Conference on, vol. 1, IEEE, 1995, pp. 588–593.
- [31] T. A. OLIVER AND R. D. MOSER, *Bayesian uncertainty quantification applied to rans turbulence models*, in Journal of Physics: Conference Series, vol. 318, IOP Publishing, 2011, p. 042032.
- [32] E. J. PARISH AND K. DURAISAMY, *A paradigm for data-driven predictive modeling using field inversion and machine learning*, Journal of Computational Physics, 305 (2016), pp. 758–774.
- [33] G. PAVLIOTIS AND A. STUART, *Multiscale methods: averaging and homogenization*, Springer Science & Business Media, 2008.
- [34] F. PEDREGOSA, G. VAROQUAUX, A. GRAMFORT, V. MICHEL, B. THIRION, O. GRISEL, M. BLONDEL, P. PRETTENHOFER, R. WEISS, V. DUBOURG, J. VANDERPLAS, A. PASSOS, D. COURNAPEAU, M. BRUCHER, M. PERROT, AND E. DUCHESNAY, *Scikit-learn: Machine learning in Python*, Journal of Machine Learning Research, 12 (2011), pp. 2825–2830.
- [35] B. PEHERSTORFER AND K. WILLCOX, *Data-driven operator inference for nonintrusive projection-based model reduction*, Computer Methods in Applied Mechanics and Engineering, 306 (2016), pp. 196–215.
- [36] M. RAISSI, P. PERDIKARIS, AND G. E. KARNADAKIS, *Physics informed deep learning (part ii): Data-driven discovery of nonlinear partial differential equations*, arXiv preprint arXiv:1711.10566, (2017).
- [37] G. ROZZA, *Reduced basis approximation and error bounds for potential flows in parametrized geometries*, Communications in Computational Physics, 9 (2011), pp. 1–48.
- [38] P. SAGAUT, *Large eddy simulation for incompressible flows: an introduction*, Springer Science & Business

- Media, 2006.
- [39] Y. SATO AND S. NAGAYA, *Evolutionary algorithms that generate recurrent neural networks for learning chaos dynamics*, in Evolutionary Computation, 1996., Proceedings of IEEE International Conference on, IEEE, 1996, pp. 144–149.
 - [40] P. J. SCHMID, *Dynamic mode decomposition of numerical and experimental data*, Journal of fluid mechanics, 656 (2010), pp. 5–28.
 - [41] G. SHULKIND, L. HORESH, AND H. AVRON, *Experimental design for non-parametric correction of misspecified dynamical models*, arXiv preprint arXiv:1705.00956, (2017).
 - [42] A. P. SINGH, S. MEDIDA, AND K. DURAIAMY, *Machine-learning-augmented predictive modeling of turbulent separated flows over airfoils*, AIAA Journal, (2017), pp. 1–13.
 - [43] S. STOLZ, N. A. ADAMS, AND L. KLEISER, *An approximate deconvolution model for large-eddy simulation with application to incompressible wall-bounded flows*, Physics of fluids, 13 (2001), pp. 997–1015.
 - [44] W. SU, M. BOGDAN, E. CANDÈS, ET AL., *False discoveries occur early on the lasso path*, The Annals of Statistics, 45 (2017), pp. 2133–2150.
 - [45] F. TAKENS ET AL., *Detecting strange attractors in turbulence*, Lecture notes in mathematics, 898 (1981), pp. 366–381.
 - [46] R. TIBSHIRANI, *Regression shrinkage and selection via the lasso*, Journal of the Royal Statistical Society. Series B (Methodological), (1996), pp. 267–288.
 - [47] R. J. TIBSHIRANI ET AL., *The lasso problem and uniqueness*, Electronic Journal of Statistics, 7 (2013), pp. 1456–1490.
 - [48] P. TREBATICĀY, *Prediction of dynamical systems by recurrent neural networks*, Inf. Sci Technol. Bull. ACHM Slovakia, 1 (2009), pp. 47–56.
 - [49] A. P. TRISCHLER AND G. M. DELEUTERIO, *Synthesis of recurrent neural networks for dynamical system simulation*, Neural Networks, 80 (2016), pp. 67–78.
 - [50] K. VEROY AND A. PATERA, *Certified real-time solution of the parametrized steady incompressible navier–stokes equations: rigorous reduced-basis a posteriori error bounds*, International Journal for Numerical Methods in Fluids, 47 (2005), pp. 773–788.
 - [51] A. WAIBEL, T. HANAZAWA, G. HINTON, K. SHIKANO, AND K. J. LANG, *Phoneme recognition using time-delay neural networks*, in Readings in speech recognition, Elsevier, 1990, pp. 393–404.
 - [52] K. WILLCOX AND J. PERAIRE, *Balanced model reduction via the proper orthogonal decomposition*, AIAA journal, 40 (2002), pp. 2323–2330.
 - [53] M. O. WILLIAMS, I. G. KEVREKIDIS, AND C. W. ROWLEY, *A data-driven approximation of the koopman operator: Extending dynamic mode decomposition*, Journal of Nonlinear Science, 25 (2015), pp. 1307–1346.
 - [54] X. XIE, M. MOHEBUJJAMAN, L. REBHOLZ, AND T. ILIESCU, *Data-driven filtered reduced order modeling of fluid flows*, arXiv preprint arXiv:1709.04362, (2017).
 - [55] J. XU AND K. DURAIAMY, *Reduced-order modeling of model rocket combustors*, in 53rd AIAA/SAE/ASEE Joint Propulsion Conference, 2017, p. 4918.