



Memory-Based Reduced Modelling and Data-Based Estimation of Opinion Spreading

Niklas Wulkow¹ · Péter Koltai¹ · Christof Schütte^{1,2}

Received: 23 June 2020 / Accepted: 12 December 2020 / Published online: 19 January 2021
© The Author(s) 2021

Abstract

We investigate opinion dynamics based on an agent-based model and are interested in predicting the evolution of the percentages of the entire agent population that share an opinion. Since these opinion percentages can be seen as an aggregated observation of the full system state, the individual opinions of each agent, we view this in the framework of the Mori–Zwanzig projection formalism. More specifically, we show how to estimate a nonlinear autoregressive model (NAR) with memory from data given by a time series of opinion percentages, and discuss its prediction capacities for various specific topologies of the agent interaction network. We demonstrate that the inclusion of memory terms significantly improves the prediction quality on examples with different network topologies.

Keywords Memory-based model · Sparse model identification · Mori–Zwanzig formalism · Nonlinear autoregressive model · Opinion dynamics · Agent-based model

Mathematics Subject Classification 37M10 · 39A50 · 91D30

1 Introduction

Political opinion polls capture how the opinions of people within a society regarding a certain topic or their current voting preferences are distributed. Individual opinions do not have to be constant, but rather are subject to change induced by impactful events or the opinions of their peers which is formalized under the term *conformity*

Communicated by Philipp M Altrock.

✉ Niklas Wulkow
niklas.wulkow@zib.de

¹ Department of Mathematics and Computer Science, Freie Universität Berlin, Berlin, Germany

² Zuse Institute Berlin, Berlin, Germany

in Stangor (2015). There have been recent advances in simulating the process in which members of a society change their opinions; see, e.g., Banisch et al. (2011), Klimek et al. (2007), Misra (2012), Li et al. (2012), Nardini et al. (2008), Böhme and Gross (2012), Bolzern et al. (2017) and the review articles Anderson and Ye (2019), Xia et al. (2011), Castellano et al. (2009), Sîrbu et al. (2017). This is in part due to increasing computing power which enables to carry out agent-based models that simulate behaviour of members of a synthetic population, such as members of a society, on the microscale by emulating the decision-making rules. The agents are often treated as the nodes of a network, while an edge between two nodes means that these agents are neighbours of each other and thus influence each other's respective opinions.

One is often not interested in modelling, or predicting, which person has which opinion, but rather, as in polls, what the percentage of each opinion within the society is. There is ample interest in deriving dynamics for the evolution of these percentages.

In this article, we will present a framework which identifies the governing equations for the dynamics of opinion percentages for different types of networks, more precisely, how the governing equations can be inferred from data on the opinion percentages. To this end, we will emulate the decision-making process with a simple agent-based model (ABM) that is based on the assumption of conformity and inspired by the ABM in Misra (2012). Introductions into agent-based modelling in general can be found in Jennings et al. (1998) and Laubenbacher et al. (2009) and specifically into agent-based models for opinion dynamics in Banisch (2016).

The literature contains a variety of approaches for finding governing equations on the macrolevel (here, opinion percentages) based on microdynamics (here, agent-based model). However, most do not deal with opinion formation or voter models, but with models originating from the context of the natural sciences. There it is well known that the aggregation process from the micro- to the macrolevel typically leads to non-Markovian processes, i.e., finding the governing equations on the macrolevel requires the inclusion of memory, cf. the Mori–Zwanzig formalism (Zwanzig 2001; Lin and Lu 2019; Chorin et al. 2002). In the context of opinion formation, this aspect is hardly discussed at all. Banisch (2014) discusses the issue for agent-based models; he gives stochastic and combinatorial arguments for the appearance of memory with heterogeneous microstructure, but does not present any practical methods for finding appropriate governing equations for the macrodynamics. Several other authors discuss the micro-macroaggregation problem in opinion formation, e.g., via influence matrices between agents (Wu et al. 2018; Ravazzi et al. 2019; De et al. 2019), but ignore memory effects entirely. Others discuss memory effects, but only on the microlevel, e.g., Jedrzejewski and Sznajd-Weron (2018), Chen et al. (2018) (agents have memory), Moussaïd et al. (2013) (agents gain experience) or Boschia et al. (2019) (microdynamics depends on collective memory). Very few articles consider the practical methods for finding governing equations on the macrolevel, e.g., by inferring them from microlevel simulation data, but memory effects are ignored, cf. Lu et al. (2019). Thus, there is a significant gap between Banisch' insight that opinion aggregation introduces memory and its practical use for finding appropriate description of the resulting macrodynamics.

This article aims at closing this gap by (1) utilizing techniques like the Mori–Zwanzig formalism and Taken’s well-known embedding theorem for showing that agent-based models for the microdynamics lead to memory effects on the macrolevel if the interaction between the agents is heterogeneous, while doing this in a way that allows for (2) proposing practical algorithmic techniques to learn governing equations for the macrodynamics including memory utilizing macroobservations of microlevel simulation data.

More precisely, we investigate complete and incomplete interaction networks: in complete networks, every agent interacts with all others (homogeneous interaction), while in incomplete networks there are subcommunities within the society that have few links between each other (heterogeneous interaction). As we will show, in the case of a complete network, one can identify a *Markovian* model for the macrodynamics of the opinion percentages using standard well-mixedness arguments known from the mean-field approaches or population limits, e.g., for predator–prey models Berryman (1992). However, arguments used for that case do not hold true in cases when the network is not complete. We will show how to use information from the past (memory) via a kind of delay embedding of the dynamics to describe the evolution of opinion percentages in the general case.

The exact reason for the inclusion of memory will formally be derived in Sect. 2 by using the Mori–Zwanzig formalism (Zwanzig 2001; Lin and Lu 2019; Chorin et al. 2002). Inspired by problems in statistical physics, the Mori–Zwanzig formalism explains how in the case of only low-dimensional observations of a high-dimensional system being available, the evolution of these observations of the full system can be obtained by replacing the missing information of the full system by past information of these available observations. This is in light of the result of Takens (1981) that states that, under fairly generic assumptions, the delay embedding of the dynamics of an observable is diffeomorphic to the dynamics of the full system.

There are various techniques for the modelling of time-discrete dynamical systems which involve the memory of the system. An intuitive approach is comprised by higher-order Markov models (Raftery 1985; Tuyen 2018). These models are defined by transition probabilities between discrete states where each state represents a sequence of cells of a discretization of the state space with a given length (“memory depth”). Although these models can be powerful in investigating the long-term behaviour of the process by means of Markov state models for Markovian processes (Bowman et al. 2014), they yield two problems: the loss of accuracy obtained from the discretization and an exponentially increasing number of states with increasing length of the sequences and number of grid cells.

Another example is simplex projection as in Sugihara and May (1990) where, using Takens’ result, subsequent states of a system are predicted from relative next steps of similar patterns as its recent history. A younger modelling technique is long short-term memory neural networks (LSTMs) (Hochreiter and Schmidhuber 1997; Pan and Duraisamy 2018) which is a subclass of recurrent neural networks and specifically designed for prediction of time series for which past information is vital. However, both these techniques provide little to no understanding of the dynamical *rules* of the system: simplex projection does not produce any model or dynamical law, but rather uses a procedure similar to the nearest neighbour classification algorithm (see, e.g., Devroye

et al. (2013)). LSTMs, as most neural networks, typically have far too many parameters to admit interpretability. An additional means for forecasting of memory-dependent dynamical systems is the well-known class of autoregressive (AR) models (Brockwell and Davis 1991), which describes the evolution of a system by a linear combination of its most recent states. Additionally, there exist variants of these AR models that are sparse (Davis et al. 2012; Fujita et al. 2007) or nonlinear (Billings 2013) or comprise both aspects in application to a singular value decomposition of a data matrix (Brunton et al. 2016). As we will see, linear (Markovian) systems cannot describe the evolution of opinion percentages even in the simplest case, but *simple* polynomial terms are sufficient for fully connected networks. We shall address this point with nonlinear AR (NAR) models, as derived through the Mori–Zwanzig formalism.

In addition to the analysis of micro-macroaggregation for opinion formation, further novelty in our work lies in the methods we propose for learning NAR models from data, to describe the evolution of opinion percentages, and their theoretical justification. We will show that the prediction accuracy of the NAR models for the opinion percentages increases with larger memory depths. To this end, we will deploy methods from data-driven (sparse) system identification—as in dynamic mode decomposition (Schmid and Sesterhenn 2008; Tu et al. 2014; Jovanovic et al. 2013) or sparse identification of nonlinear dynamics (SINDy) (Brunton et al. 2016a)—to the field of opinion dynamics. More precisely, we will extend SINDy towards finding (sparse) NAR models to describe the evolution of opinion percentages. The new method is called “sparse identification of nonlinear autoregressive models” (SINAR), as it is technically a natural generalization of SINDy by including nonlinear memory terms. We will demonstrate that SINAR is well suited for our purposes in learning macroscopic opinion dynamics. A conceptually similar method has been introduced in Brunton et al. (2016) with Hankel alternative view of Koopman (HAVOK). It can be interpreted as a special case of SINAR.

Outline In Sect. 2, we start with outlining the opinion aggregation process and proceed with the derivation of NAR models for the evolution of observations through the Mori–Zwanzig formalism. Next, in Sect. 3, we present the SINAR method for estimating the coefficients in these NAR models from data. Last, we demonstrate how to apply SINAR for increasing the accuracy of prediction of opinion percentages in the case of incomplete interaction networks in Sect. 4.

2 Derivation of a Nonlinear Autoregressive Model Using the Mori–Zwanzig Formalism

Below, we will model the spread of opinions inside a closed society by an agent-based model. It will consist of a high number N of agents who change their opinions X_i , $i = 1, \dots, N$, within a finite set of M possible opinions over discrete time steps according to a rule that is based on the opinions of themselves and other agents. This rule will be Markovian, or memory-free, i.e., the changes of opinions are only influenced by opinions in the current time step. These dynamics will be called the *microdynamics*. The state of the microdynamics at time t is denoted by $X_t = [(X_t)_1, \dots, (X_t)_N]^T$. The respective state space is denoted by \mathbb{X} and has cardinality $|\mathbb{X}| = M^N$.

We will only be able to observe the percentages of opinions, i.e., the ratios of those among all agents with each of the M opinions. In this article, we are interested in identifying the dynamical rules of the evolution of the percentages of opinions, which we call the *macrodynamics*. Identifying the dynamics of low-dimensional observations of a higher-dimensional system is a typical setup for the Mori–Zwanzig formalism (Zwanzig 2001; Chorin et al. 2002; Lin and Lu 2019). We will consider a general framework for this and show how it yields a nonlinear autoregressive model (Billings 2013) for the macrodynamics. Later on, we show how it can be applied to the specific case of the spread of opinions.

2.1 The Setting: Microdynamics and Projected Observations

First we assume that the microdynamics are Markovian (memory-free) and deterministic. We consider the dynamical system $F : \mathbb{X} \rightarrow \mathbb{X}$ that governs the microdynamics

$$X_{t+1} = F(X_t) \in \mathbb{X}. \tag{2.1}$$

Further, we denote the space of observations of the microdynamics (observables) by $\mathbb{Y} \subseteq \mathbb{R}^m$ and by $\mathcal{G} := \{g : \mathbb{X} \rightarrow \mathbb{Y}\}$ the set of functions that map states of the dynamical system (2.1) to \mathbb{Y} . We suppose from here on that we do not have knowledge of the state of the microdynamics at any point in time, but instead only have the value of the fixed observable $x = \xi(X) \in \mathbb{Y}$ which we call the accessible, or *relevant*, variables.

Additionally, we define the subspace \mathcal{H} of functions in \mathcal{G} that depend only on these relevant variables and map to \mathbb{Y} as $\mathcal{H} := \{h \in \mathcal{G} \mid \exists \tilde{h} : \xi(\mathbb{X}) \rightarrow \mathbb{Y} : h = \tilde{h} \circ \xi\}$. Functions in \mathcal{H} still depend on $X \in \mathbb{X}$, but the information of $\xi(X)$ is enough to evaluate them. When we write $h(x)$ for $x \in \mathbb{Y}$, we abuse notation and mean $h(\xi(x))$. An example is

$$\mathbb{X} = \mathbb{R}^2, \quad \xi(X) = X_1 + X_2, \quad h(X_1, X_2) = (X_1 + X_2)^2 = \xi(X)^2.$$

In this case, it is enough to know the value of $\xi(X)$ to evaluate $h(X)$.

The goal is now to represent the evolution of the observations $x_t = \xi(X_t)$ under the microdynamics with knowledge only about values of x_t , but not of the states X_t of the microdynamics. As illustrated in the following diagram, instead of taking one step of the microdynamics and then evaluating ξ , we only have access to the observation $\xi(X)$ and want to evaluate $\xi(F(X))$ under the premise that $\xi(X) = x$.

$$\begin{array}{ccc}
 X & \xrightarrow{F} & F(X) \\
 \xi \downarrow & & \downarrow \xi \\
 \xi(X) = x & \xrightarrow{?} & \xi(F(X))
 \end{array} \tag{2.2}$$

To this end, we define a projection operator $P : \mathcal{G} \rightarrow \mathcal{H}$ that maps a function depending on X to a function depending on $\xi(X)$. We additionally define its complement $Q :=$

$Id - P$. We assume from now on that the microdynamics are stationary with an F -invariant probability distribution μ over \mathbb{X} , so that when asking what $g(X)$ is, we assume that X_t is distributed by μ .¹ We, of course, are interested in the case $g = \xi \circ F$. We follow Lin and Lu (2019) until the end of Sect. 2.2 and define P as the orthogonal projection onto the span of a set of linearly independent functions from \mathcal{H} . These functions are denoted by $\varphi_1, \dots, \varphi_L : \mathbb{Y} \rightarrow \mathbb{R}^m$ which build the columns of $\varphi = [\varphi_1, \dots, \varphi_L]$.

$$(Pg)(x) := \varphi(x)\langle\varphi, \varphi\rangle^{-1}\langle\varphi, g\rangle \tag{2.3}$$

where $x \in \mathbb{Y}$ and the scalar product $\langle \cdot, \cdot \rangle$ is defined for matrix-valued functions $f : \mathbb{X} \rightarrow \mathbb{R}^{m \times a}$ and $g : \mathbb{X} \rightarrow \mathbb{R}^{m \times b}$ as

$$\langle f, g \rangle := \int_{\mathbb{X}} \underbrace{f(X)^T}_{\in \mathbb{R}^{a \times m}} \underbrace{g(X)}_{\in \mathbb{R}^{m \times b}} d\mu(X) \in \mathbb{R}^{a \times b},$$

which itself is matrix-valued. The term $\langle\varphi, \varphi\rangle$ is a mass matrix that ensures that P is an orthogonal projection. This orthogonal projection has the property that Pg is the closest function in $span(\varphi)$ to g with respect to $\langle \cdot, \cdot \rangle$.

Note that if \mathcal{H} is infinite-dimensional, one would need an infinite number of functions to yield that $span(\varphi) = \mathcal{H}$. In this case, the projection formalism is well defined if \mathcal{H} is closed. In practice, in this case for the computation that will follow one would choose a sufficiently rich finite set of functions so that $span(\varphi)$ covers those parts of \mathcal{H} that are of interest.

2.2 Mori–Zwanzig Representation of the Macrodynamics

We will now show how to represent the evolution of the observations over time. With the Koopman operator (Koopman 1931) \mathcal{K} for the system (2.1), defined as the operator that maps a function $g \in \mathcal{G}$ to $g \circ F \in \mathcal{G}$, we consider the Dyson formula

$$\mathcal{K}^{t+1} = \sum_{k=0}^t \mathcal{K}^{t-k} P\mathcal{K}(Q\mathcal{K})^k + (Q\mathcal{K})^{k+1}. \tag{2.4}$$

The Dyson formula describes a way to iteratively split up the application of the Koopman operator to a function g into parts $P\mathcal{K}g$ and $Q\mathcal{K}g$. Equation (2.4) yields, by application of both sides of the equation to ξ and evaluation at the initial value X_0 of the microdynamics, that

$$x_{t+1} = \sum_{k=0}^t [P(\rho^k \circ F)](x_{t-k}) + \rho^{t+1}(X_0). \tag{2.5}$$

¹ A natural candidate for P would be the conditional expectation with respect to μ , given by $(Pg)(x) = \mathbb{E}[g(X) \mid \xi(X) = x]$; see Appendix A.4. Approximating the conditional expectation can be a challenging task, see Gilani et al. (2020). Instead, we consider the orthogonal projection onto basis functions since we are seeking models spanned by such functions with the option to control the sparsity of the model. In Chorin et al. (2002), the connection between both projections is discussed.

where $\rho^k := (Q\mathcal{K})^k \xi$. The derivation of Eq. (2.5) is explained in detail in Appendix A.1, together with interpretation of terms of its right-hand side.

Substituting the definition of P as the orthogonal projection onto basis functions as in (2.3), we obtain

$$P(\rho^k \circ F)(x_{t-k}) = \varphi(x_{t-k}) \langle \varphi, \varphi \rangle^{-1} \langle \varphi, \rho^k \circ F \rangle =: \varphi(x_{t-k}) h_k \in \mathbb{R}^m \tag{2.6}$$

with vector-valued coefficients $h_k = \langle \varphi, \varphi \rangle^{-1} \int_{\mathbb{X}} \varphi(\xi(X))^T \rho^k(F(X)) d\mu(X)$.

Finding a suitable approximation of the non-accessible noise term $\rho^{t+1}(X_0)$ in (2.5) is generally a non-trivial task and depends on properties of the microdynamics. Examples are discussed in Li and Chu (2017), Hijón et al. (2010), Kondrashov et al. (2015). From this point onwards, we will make the simplification of replacing $\rho^{t+1}(X_0)$ by a zero-mean stochastic noise term $\varepsilon_{t+1} \in \mathbb{R}^m$. A typical practice is to let ε_{t+1} be a zero-mean Gaussian random variable as, e.g., in Lin and Lu (2019), Lei et al. (2016). With this, we obtain the macrodynamics

$$x_{t+1} = \sum_{k=0}^t \varphi(x_{t-k}) h_k + \varepsilon_{t+1}. \tag{2.7}$$

As we can see, the evolution of the observations now depends on past terms, although the microdynamics are Markovian. For $k > 0$, the terms $[P(\rho^k \circ F)](x_{t-k})$ in Eq. (2.5) and $\varphi(x_{t-k})$ in Eq. (2.7) are usually referred to as *memory terms*.

2.3 Macrodynamics as a Nonlinear Autoregressive Process

If it is reasonable to assume a sufficiently fast decay of the terms h_k with increasing k , the memory terms that lie far in the past have negligible influence (Horenko et al. 2007; Venkataramani et al. 2017; Chorin et al. 2000; Zhu et al. 2018). In light of (2.5) and (2.6), it is sufficient that the ρ^k decay fast. To understand when this is the case, we recall $\rho^k = (Q\mathcal{K})^k \xi$ and assume the range(P) $\approx \mathcal{H}$, i.e., functions parametrized by ξ are well approximated by the chosen approximation space. Then, ρ^k decays fast if $Q\mathcal{K}$ has a small norm, which is the case if F mixes well functions that are perpendicular to \mathcal{H} . In other words, the dominant modes of \mathcal{K} should align well with the space \mathcal{H} . For quantitative statements we refer to Zhu et al. (2018).

Thus, in order to obtain a feasible number of memory terms, from now on we approximate the dynamics by ending the sum in (2.7) with $k = p - 1$ instead of $k = t$, i.e., by truncating the terms $\varphi(x_{t-p})h_p, \dots, \varphi(x_0)h_t$. Regarding the selection of an appropriate value for the *memory depth* p , there are various methods such as Information Criteria (Konishi and Kitagawa 2008; Aho et al. 2014) or the L-curve method (Hansen and D. O’leary 1993). We have thus derived a nonlinear autoregressive model (NAR) (Billings 2013; An and Huang 1996) over x given by

$$x_{t+1} = \sum_{k=0}^{p-1} \varphi(x_{t-k}) h_k + \varepsilon_{t+1} \tag{2.8}$$

with matrix-valued basis functions and vector-valued coefficients h_k .

In Sect. 3, we will introduce a method that identifies coefficients for NAR models in a way that is motivated by system identification methods such as dynamic mode decomposition (Williams et al. 2014; Tu et al. 2014), extended dynamic mode decomposition (Williams et al. 2014) or sparse identification of nonlinear dynamics (Brunton et al. 2016a, b), see Fig. 1, where the dynamics are expressed with a vector of scalar-valued basis functions and a matrix-valued coefficient. Having selected the scalar-valued basis functions $\tilde{\varphi}_1, \dots, \tilde{\varphi}_K$ and denoting $\tilde{\varphi} = [\tilde{\varphi}_1, \dots, \tilde{\varphi}_K]^T : \mathbb{Y} \rightarrow \mathbb{R}^K$, we thus formulate the macrodynamics

$$x_{t+1} = \sum_{k=0}^{p-1} H_k \tilde{\varphi}(x_{t-k}) + \varepsilon_{t+1}, \quad (2.9)$$

with $H_k \in \mathbb{R}^{m \times K}$. Although seeming like only a slight notational modification, both formulations represent different model forms. While in (2.8) the dynamics are expressed using different basis functions and the same coefficients across all coordinates, we will now switch to the framework in (2.9) where we select scalar-valued basis functions $\tilde{\varphi}_1, \dots, \tilde{\varphi}_L$ which are used for each coordinate, while the coefficients for all coordinates can be different (the different rows of the H_k). In summary, for (2.8), one chooses L m -dimensional basis functions and finds L -dimensional coefficients, while for (2.9), one chooses K one-dimensional basis functions and finds $(m \times K)$ -dimensional coefficients.

Equation (2.9) is still consistent with the way we derive (2.8) through the Mori–Zwanzig formalism: basis functions are evaluated at observations made at distinct times—no terms with mixed delays occur. In Appendix A.2, we show how to choose basis functions and coefficients in each of the models to derive the equivalent dynamics. Please note that this does not mean that both model forms are always equivalent, as explained above. Merely, one can always choose $\tilde{\varphi}$ in dependence on φ , respectively, vice versa, in a way that makes the dynamics equivalent.

2.4 Stochastic Microdynamics

Let us consider stochastic dynamics

$$X_{t+1} = F(X_t, \omega_t)$$

where $\omega_t \in \Omega$ is a random influence on F which is now defined as $F : \mathbb{X} \times \Omega \rightarrow \mathbb{X}$. We will assume that the noise process ω_t , $t \in \mathbb{N}$, is i.i.d. with law \mathbb{P} . In this case, we only strive to forecast the *expected* macrodynamics, and define the (stochastic) Koopman operator as

$$(\mathcal{K} \circ g)(X) = \mathbb{E}_{\mathbb{P}}[g(F(X, \omega))].$$

The spaces \mathcal{G} and \mathcal{H} , just as the projection P remain unchanged. Naturally, to the derivation of the Mori–Zwanzig approximation we need to apply the necessary obvious modifications. For example, the last step in (2.5) now has to be modified as:

$$\left[PK\rho^k \right] (x_{t-k}) = \varphi(x_{t-k}) \langle \varphi, \varphi \rangle^{-1} \int_{\Omega} \int_X \varphi(\xi(X))^T \rho^k(F(X, \omega)) d\mu(X) d\mathbb{P}(\omega).$$

We can thus obtain the identical structure of the macrodynamics as in (2.7) where for the computation of the coefficients h_k in (2.6) the expectation with respect to \mathbb{P} had to be added.

3 Sparse Identification of Nonlinear Autoregressive Models (SINAR)

We propose here a method of data-based identification for coefficients H_k in (2.7) that is an extension of the sparse identification of nonlinear dynamics (SINDy) algorithm from Brunton et al. (2016a), Brunton et al. (2016b), Kaiser et al. (2018). SINDy can be used to identify the governing equations of a Markovian—in our case, discrete time—dynamical system

$$x_{t+1} = f(x_t) \in \mathbb{R}^m \tag{3.1}$$

from data

$$\mathbf{X} = \begin{bmatrix} | & & | \\ x_0 & \dots & x_{T-1} \\ | & & | \end{bmatrix}, \mathbf{X}' = \begin{bmatrix} | & & | \\ x_1 & \dots & x_T \\ | & & | \end{bmatrix}, \mathbf{X}, \mathbf{X}' \in \mathbb{R}^{m \times T}.$$

We will extend this method to non-Markovian systems by applying SINDy to an extended version of \mathbf{X} , the *Hankel* matrix

$$\tilde{\mathbf{X}} = \begin{bmatrix} x_{p-1} & \dots & x_{T-1} \\ \vdots & & \vdots \\ x_0 & \dots & x_{T-p} \end{bmatrix}.$$

In essence, this is the concept used for the Hankel alternative view of Koopman (HAVOK) analysis from Brunton et al. (2016), where an autoregressive model is identified on transformed coordinates obtained from a singular value decomposition of the Hankel matrix from a scalar-valued observation function to separate linear from nonlinear, or even chaotic, behaviour of a Markovian system. We, however, seek a formulation for the dynamics of multidimensional observations. In this section and by the choice of the name SINAR, we explicitly want to point out the connection of system identification methods for nonlinear Markovian systems to their counterparts for nonlinear non-Markovian systems (with finite memory these are NAR systems) that can be derived through the Mori–Zwanzig formalism from Sect. 2.

3.1 SINDy: A Short Summary

We start with a short description of SINDy (Brunton et al. 2016a). In SINDy, we try to approximate each coordinate of f by a linear combination of basis functions $\theta_i : \mathbb{R}^m \rightarrow \mathbb{R}$ and define

$$\Theta(x) = \begin{bmatrix} \theta_1(x) \\ \vdots \\ \theta_v(x) \end{bmatrix}, \quad \Theta(\mathbf{X}) = \begin{bmatrix} \theta_1(x_0) \dots \theta_1(x_{T-1}) \\ \vdots \\ \theta_v(x_0) \dots \theta_v(x_{T-1}) \end{bmatrix}.$$

To this end, we fit a sparse coefficient matrix $\Xi \in \mathbb{R}^{m \times v}$ with rows Ξ_i to the data \mathbf{X} , \mathbf{X}' by solving for every row \mathbf{X}'_i of \mathbf{X}' ,

$$\Xi_i = \arg \min_{\Xi_i} \|\mathbf{X}'_i - \Xi_i \Theta(\mathbf{X})\|_F + \lambda \|\Xi_i\|_1. \quad (3.2)$$

We then obtain the model

$$x_{t+1} \approx \Xi \Theta(x_t). \quad (3.3)$$

In (3.2), we enforce a sparsity constraint using the LASSO regression algorithm (Tibshirani 1996) in which a regularization term is added onto the coefficient matrix, in order to only obtain the basis functions from Θ that are dominant for the relation between x_{t+1} and $\Theta(x_t)$.

The use of the 1-norm generates a sparse solution if we set $\lambda > 0$ appropriately. Sparse models will often times be less accurate than non-sparse models. However, what we gain through a sparse right-hand side of (3.3) is a better interpretability of the model since only the dominant terms have been identified as influential to the dynamics. It is vital to set λ so that the loss of accuracy is minimal compared to the gain in interpretability.

SINDy is closely related to the (first step of) the method of dynamic mode decomposition (DMD) (Williams et al. 2014; Tu et al. 2014), which aims at finding a linear connection between x_t and x_{t+1} . To this end, one solves²

$$A = \arg \min_A \|\mathbf{X}' - A\mathbf{X}\|_F. \quad (3.4)$$

3.2 Extending SINDy to SINAR

When the dynamical model (3.1) is insufficient in the sense that x_{t+1} depends not only on x_t but on memory terms too, we can apply the SINDy algorithm to suitably transformed data to obtain a nonlinear autoregressive model as in (2.9) with sparse coefficients. That is, only a few basis functions should occur with nonzero coefficients.

² In a second step, DMD then uses Ξ from (3.2) to uncover properties of the Koopman operator of the system. SINDy, instead, tries to explain the evolution of x_t by basis functions that do not have to be linear. Still, essentially, the problem (3.4) is equivalent to (3.2) for $\Theta(x) = x$ and $\lambda = 0$. Further, there exists a sparse version of DMD (Jovanovic et al. 2013), where the sparsity constraint is enforced by the additive 1-norm regularization as in (3.2). Then the emerging minimization problem is the same as (3.2) with $\Theta(x) = x$.

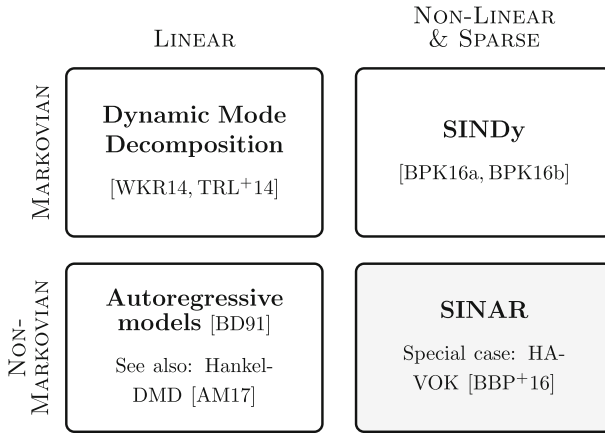


Fig. 1 Relation between different system identification methods. All of them are based on solving a least squares problem with respect to transformations of past to future states. While the AR minimization problem can be seen as the DMD problem on delay-embedded states and SINDy finds a nonlinear instead of linear connection between states (as in Hankel-DMD in Arbabi and Mezic (2017)), SINAR finds a nonlinear connection between multiple past states and future ones. SINAR allows for imposing a sparsity constraint onto the determination of macromodels in the same fashion as is done in SINDy for Markovian systems. This has already been done in a special way in Brunton et al. (2016), which is a special case of SINAR

Selecting a memory depth p and denoting

$$\tilde{x}_t := \begin{bmatrix} x_t \\ \vdots \\ x_{t-p+1} \end{bmatrix} \in \mathbb{R}^{mp},$$

let us define as data matrices the Hankel matrix

$$\tilde{\mathbf{X}} = \begin{bmatrix} x_{p-1} & \dots & x_{T-1} \\ \vdots & & \vdots \\ x_0 & \dots & x_{T-p} \end{bmatrix} = \begin{bmatrix} | & & | \\ \tilde{x}_{p-1} & \dots & \tilde{x}_{T-1} \\ | & & | \end{bmatrix} \in \mathbb{R}^{mp \times (T-p+1)} \tag{3.5}$$

and $\mathbf{X}' = \begin{bmatrix} | & & | \\ x_p & \dots & x_T \\ | & & | \end{bmatrix} \in \mathbb{R}^{m \times (T-p+1)}.$

Again, we choose basis functions

$$\tilde{\Theta}(\tilde{x}) = \begin{bmatrix} \tilde{\theta}_1(\tilde{x}) \\ \vdots \\ \tilde{\theta}_v(\tilde{x}) \end{bmatrix}, \quad \tilde{\Theta}(\tilde{\mathbf{X}}) = \begin{bmatrix} \tilde{\theta}_1(\tilde{x}_{p-1}) & \dots & \tilde{\theta}_1(\tilde{x}_{T-1}) \\ \vdots & & \vdots \\ \tilde{\theta}_v(\tilde{x}_{p-1}) & \dots & \tilde{\theta}_v(\tilde{x}_{T-1}) \end{bmatrix}$$

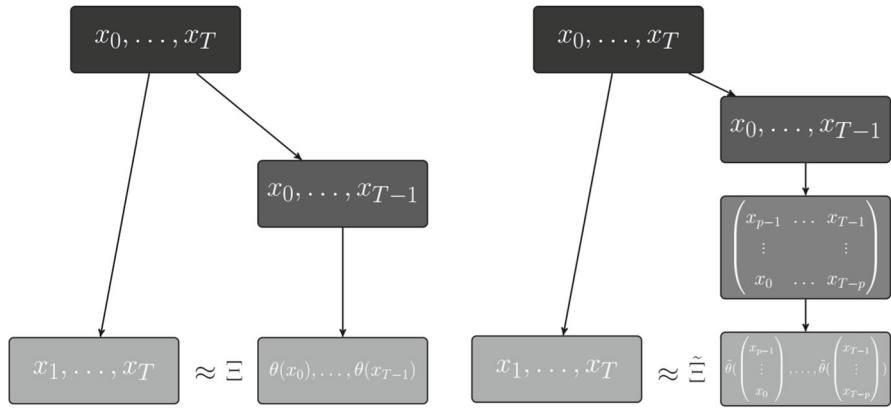


Fig. 2 Sketch of the SINDy algorithm (left) and SINAR (right). SINAR contains the additional step of creating a Hankel matrix

for example

$$\tilde{\Theta}(\tilde{x}_t) = [(x_t)_1^2, (x_t)_1(x_t)_2, \dots, \sin((x_{t-1})_1), \dots, (x_{t-2})_m(x_{t-3})_1]^T,$$

and minimize for every row $\tilde{\Xi}_i$ of $\tilde{\Xi}$:

$$\tilde{\Xi}_i = \arg \min_{\tilde{\Xi}_i} \|\mathbf{X}'_i - \tilde{\Xi}_i \tilde{\Theta}(\tilde{\mathbf{X}})\|_F + \lambda \|\tilde{\Xi}_i\|_1. \tag{3.6}$$

Then with the basis functions with nonzero coefficients in $\tilde{\Xi} \in \mathbb{R}^{m \times v}$, we have derived a nonlinear autoregressive model that approximates the evolution of x :

$$x_{t+1} = \tilde{\Xi} \tilde{\Theta}(\tilde{x}_t) \in \mathbb{R}^m, \quad \text{or, equivalently,} \quad (x_{t+1})_i = \sum_{j=1}^v \tilde{\Xi}_{ij} \tilde{\theta}_j(\tilde{x}_t). \tag{3.7}$$

By deleting all columns of $\tilde{\Xi}$ that only contain zeros, which should be many if we enforce the sparsity constraint, we get a reduced matrix and thus a low number of terms on the right-hand side of (3.7). We have thus identified a sparse nonlinear autoregressive model so that we call this extension of SINDy sparse identification of nonlinear autoregressive models (SINAR). Note that for a memory depth of $p = 1$, SINDy and SINAR are equivalent. Figure 1 shows the connections between several prominent methods for learning macrodynamics from microsimulation data in the Markovian and non-Markovian setting. Figure 2 further illustrates the different structures of SINDy and SINAR.

The choice of $\tilde{\Theta}$ allows for an arbitrary functional dependence between the distinct time-delayed observables. We can recover the special structure used in the Mori–Zwanzig formalism (2.8) and (2.9) by a particular choice of the basis by choosing

$$\tilde{\Theta}(\tilde{x}_t) = [\tilde{\varphi}_1(x_t), \dots, \tilde{\varphi}_K(x_t), \dots, \tilde{\varphi}_1(x_{t-p+1}), \dots, \tilde{\varphi}_K(x_{t-p+1})]^T,$$

with $\tilde{\varphi}_1, \dots, \tilde{\varphi}_K$ being scalar-valued functions as introduced in Sect. 2.3. Then we could directly estimate the coefficients H_k of the model (2.9)—which was derived through the Mori–Zwanzig formalism previously—from data, provided its distribution is approximately μ . Then $\tilde{\Xi}$ has the block-wise form

$$\tilde{\Xi} = [H_0, \dots, H_{p-1}] \in \mathbb{R}^{m \times pK}$$

and

$$\tilde{\Xi} \tilde{\Theta}(\tilde{x}_t) = \sum_{k=0}^{p-1} H_k \begin{bmatrix} \tilde{\varphi}_1(x_{t-k}) \\ \vdots \\ \tilde{\varphi}_K(x_{t-k}) \end{bmatrix}.$$

Of course, by choosing linear basis functions $\tilde{\Theta}(\tilde{x}_t) = \tilde{x}_t$ and setting $\lambda = 0$, one obtains a well-known linear autoregressive model (Brockwell and Davis 1991). Except for the sparsity term, the determination of model coefficients as in (3.6) is exactly the least squares method commonly used for the linear AR models. In Appendix A.3, we explain the structural equivalences and differences between SINDy, SINAR, DMD and AR models that are also sketched in Fig. 1.

The covariance of the noise term ε_{t+1} in (2.9) can be estimated in the common way for linear or nonlinear AR models (Brockwell and Davis 1991; Lin and Lu 2019) by calculating the statistical covariance between \mathbf{X}' and $\tilde{\Xi} \tilde{\Theta}(\mathbf{X})$ (see Appendix A.6 for more details on both statements).

In Appendix B, we apply SINAR to an extended Hénon system, a two-dimensional dynamical system that admits a global attractor, and inspect both its accuracy in short-term predictions and its capacity to reconstruct the original attractor. This is to illustrate basic properties of nonlinear autoregressive models for a simple system yielding complex dynamics.

4 Application to an Agent-Based Model for Opinion Dynamics

We will now consider a network-based model of agents that change their opinions on a topic based on the opinions of their neighbours in the network. Suppose, we can only observe the percentages of agents inside the network that share each opinion, but not which agent exactly has which opinion, as in an anonymous opinion poll. Describing the evolution of these percentages can be approached by the Mori–Zwanzig formalism that we discussed in Sect. 2, since they are simply observations of hidden microdynamics. We will demonstrate the efficacy of NAR models in predicting the evolution of opinion percentages, compared with Markovian models. We use a time-discrete agent-based model (ABM), similar to the concept of modelling opinion changes in a population explained in Misra (2012). The ABM in Misra (2012), however, is time-continuous, while we use a time-discretized version of it. To apply the Mori–Zwanzig formalism to a time-continuous microdynamics, we refer the interested reader to the literature such as (Chorin et al. 2000, 2002).

4.1 Formulating the ABM

The ABM is given as follows: suppose there are N agents and each agent has exactly one out of M different opinions, denoted by $1, \dots, M$. The vector X_t , which comes from

$$\mathbb{X} = \{1, \dots, M\}^N,$$

then represents the opinions of each agent at time t and $(X_t)_i$ denotes the opinion of agent i at time t . The neighbourhoods of all agents are represented by the symmetric adjacency matrix $A \in \{0, 1\}^{N \times N}$ where $A_{ij} = 1$ means that agents i and j are neighbours of each other and $A_{ij} = 0$ otherwise. Let $N_i := \#\{j : A_{ij} = 1\}$ be the number of neighbours of an agent. The diagonal entries of A are set to 1, so that every agent is its own neighbour.

Let the procedure of opinion changing be given by the following rule: in every time step, every agent picks one of its neighbours in the network uniformly at random and changes its opinion with *adaption probability* $\alpha_{m'm''}$ where m' is the opinion of the agent and m'' is the opinion of the selected neighbour. This results in the term

$$\mathbb{P}[(X_{t+1})_i = m'' | (X_t)_i = m'] = \alpha_{m'm''} \frac{\#\{j : A_{ij} = 1 \text{ and } (X_t)_j = m''\}}{N_i} \text{ for } m' \neq m'',$$

which we denote by $p_i^t(m', m'')$. The probability for an agent not to change its opinion thus is

$$p_i^t(m', m') = \mathbb{P}[(X_{t+1})_i = m' | (X_t)_i = m'] = 1 - \sum_{m'' \neq m'} p_i^t(m', m'').$$

In algorithmic form, the agent-based model is executed in the following way:

Algorithm 1: Agent-based opinion change model

- 1 Choose end time T , number of agents N , network adjacency matrix A , opinion change coefficients $\alpha_{m'm''}$, initial opinions X_0
 - 2 **for** $t = 0, \dots, T$ **do**
 - 3 **for** $i = 1, \dots, N$ **do**
 - 4 Draw j from $\{j : A_{ij} = 1\}$ uniformly at random (Choose neighbour)
 - 5 Draw $u_i \sim \mathcal{U}[0, 1]$
 - 6 If $u_i < \alpha_{(X_t)_i(X_t)_j} : (X_{t+1})_i = (X_t)_j$ (Adapt neighbour's opinion)
 - 7 **end**
 - 8 **end**
-

To clarify the notation, remember that $(X_t)_i$ and $(X_t)_j$ denote the opinions of agents i and j at time t . Hence, $\alpha_{(X_t)_i(X_t)_j}$ is the adaption probability of opinion $(X_t)_j$ given that an agent has opinion $(X_t)_i$. Note that in each time t every agent is given the opportunity to change its opinion, and whether this happens is a probabilistic event depending only on the opinions at time t .

We can now state the so-defined microdynamics by

$$X_{t+1} = F(X_t, \omega_t)$$

where at every time step, ω_t denotes a tuple consisting of N agents that represents the chosen neighbour of each agent plus numbers $u_i \sim \mathcal{U}[0, 1]$ that govern the adaption probability $\alpha_{(X_t)_i, (X_t)_j}$ as in Algorithm 1. To be more precise, ω_t has the form

$$\omega_t = [j_1, \dots, j_N, u_1, \dots, u_N], \quad j_i \sim \mathcal{U}\{j : A_{ij} = 1\}, \quad u_i \sim \mathcal{U}[0, 1].$$

F then is given by

$$(X_{t+1})_i = F(X_t, \omega_t)_i = \begin{cases} (X_t)_{j_i} & \text{if } u_i < \alpha_{(X_t)_i, (X_t)_{j_i}} \\ (X_t)_i & \text{otherwise.} \end{cases}$$

This way of stating the microdynamics seems complicated compared to the more intuitive option of denoting by $(\omega_t)_i$ the new opinion of the i th agent, distributed by $[p_i^t((X_t)_i, 1), \dots, p_i^t((X_t)_i, M)]$. However, this would mean that the distribution of ω_t changes over time, since the p_i^t depend on $(X_t)_i$. For the Mori–Zwanzig formalism, this would prevent us from applying the procedure of skew–shift systems introduced in Sect. 2.4 where we drew all ω_t a priori and thus independently of the X_t . By using the notation of ω_t denoting a tuple of neighbours j_i and random numbers u_i that are compared to the adaption coefficients, we can draw the whole sequence of ω_t independently of the X_t and maintain consistency with the notation of skew–shift systems.

4.2 Deducing Macrodynamics from the ABM

Closed-form macrodynamics.

We now define as the *opinion percentages* the function

$$\xi(X) = \frac{1}{N} \begin{bmatrix} \#X_i = 1 \\ \vdots \\ \#X_i = M \end{bmatrix}$$

and are interested in modelling how these percentages evolve over time. It turns out that for a complete network, i.e., $A_{ij} = 1 \forall i, j$, we can derive macrodynamics for the expected evolution of

$$x_t := \xi(X_t),$$

that do not require memory terms. They are given by

$$\mathbb{E}[(x_{t+1})_{m'} \mid x_t] = (x_t)_{m'} + \sum_{m'' \neq m'} (\alpha_{m''m'} - \alpha_{m'm''})(x_t)_{m''} (x_t)_{m'} \text{ for } m' = 1, \dots, m. \tag{4.1}$$

This equation can be derived as follows: in case of a complete network, $p_i^t(m', m'') \equiv p^t(m', m'')$ is independent of i because the percentages of opinions among neighbours are equal for all agents since they all have the same neighbours. Then

$$p^t(m', m'') = \alpha_{m'm''}(x_t)_{m''}.$$

In every time step, every agent with opinion m' chooses its opinion in the next time step with respective probabilities $p^t(m', m'')$ for all opinions $m'' \neq m'$ and probability $1 - \sum_{m'' \neq m'} p^t(m', m'')$ for keeping opinion m' . Since the number of these agents is given by $N \cdot (x_t)_{m'}$, the expected absolute number of agents that change their opinion from m' to m'' is given by

$$\begin{aligned} & \mathbb{E}[\#\text{Agents changing opinion from } m' \text{ to } m''] \\ &= \sum_{i:(X_t)_i=m'} p^t(m', m'') \\ &= N \cdot (x_t)_{m'} \cdot p^t(m', m'') \\ &= N \cdot (x_t)_{m'} \cdot \alpha_{m'm''} \cdot (x_t)_{m''}. \end{aligned}$$

This is the expected absolute number of agents that change their opinion from m' to m'' . This means that from this term alone, the percentage $(x_t)_{m'}$ of m' is reduced by $\frac{1}{N}$ times this term, which is $\alpha_{m'm''}(x_t)_{m''}(x_t)_{m'}$. Since at the same time agents with opinion m'' can change their opinion to m' with probability $\alpha_{m''m'}(x_t)_{m'}(x_t)_{m''}$, we have to subtract the analogous term for $\mathbb{E}[\#\text{Agents changing opinion from } m'' \text{ to } m']$ and the factor $(\alpha_{m''m'} - \alpha_{m'm''})$ comes in. As a consequence, for a complete network the expected evolution of x can be written in terms of x alone, without requiring additional information of the microstate X .

Consequences of the Mori–Zwanzig formalism.

In the abstract language of the Mori–Zwanzig formalism from Sect. 2, the above means that

$$PK\xi = K\xi, \quad \text{thus } QK\xi = 0, \tag{4.2}$$

because we can express $K\xi = \mathbb{E}[\xi \circ F]$ as a function of ξ directly by using (4.1). Let us now consider (2.5), where terms of the form

$$PK\rho^k \quad \text{with } \rho^k = (QK)^k \xi$$

occur. Equation (4.2) yields for $k > 0$ that $\rho^k = (QK)^{k-1}(QK\xi) = 0$. In this way, we can see that memory terms are not required for the dynamics of ξ if the network is complete. However, this is generally not the case for incomplete networks, as demonstrated in detail in Banisch (2014). In other words, (4.2) is no longer valid so that the ρ^k do not vanish. In this case, by using as P the orthogonal projection onto basis functions we were able to find approximate representations of the terms $P(\rho^k \circ F)$ in (2.5). Here lies another part of the value of the application of the Mori–Zwanzig formalism: it installs that the structure of the ensuing macrodynamics in (2.5) is additive, i.e., it can be written as a sum of transformations of memory terms of *individual*

delays, as opposed to memory terms containing mixed delays (e.g., $\psi_1(x_t)\psi_2(x_{t-1})$). This guides our choice for a good approximation structure and reduces the number of potential basis functions from exponential in the delay depth p to linear.³

For an incomplete network which is still sufficiently densely connected, we expect the microdynamics to be in expectation still close to that of a complete network. Thus, in such a case we expect $Q\mathcal{K}\xi \approx 0$, even if (4.2) does not hold exactly. Consequently, assuming dense connectedness, the opinion percentages should allow for a closed-form description of their evolution with a small memory depth. In the following, we will use SINAR to identify NAR models of this form suggested by the Mori–Zwanzig formalism.

4.3 Recovering the Macrodynamics in Case of an Incomplete Network

We now create realizations of the ABM with networks that consist of equally sized clusters of agents. Edges between agents from different clusters exist, but are few. Inside the clusters, all agents are connected with each other. To this end, we create networks with a total number of agents N consisting of equally sized clusters. Two agents from different clusters are connected with probability $p_{between}$.

From the same initial state and with the same parameters, we create multiple realizations of the form $[X_0 \dots, X_T]$ of the ABM and deduce the percentages of opinions $[x_0, \dots, x_T] = [\xi(X_0), \dots, \xi(X_T)]$. We denote the realizations of the resulting macrodynamics by $\mathbf{X}_1, \dots, \mathbf{X}_r$ and divide these data into training data $\mathbf{X}_1, \dots, \mathbf{X}_{train}$ and validation data $\mathbf{X}_{train+1}, \dots, \mathbf{X}_r$. Subsequently, we execute the SINAR method with different memory depths p on the training data. SINAR gives us NAR models that we use for the reconstruction of the validation data. For this, the SINAR method can straightforwardly be modified for multiple trajectories by defining data matrices $\mathbf{X}' = [\mathbf{X}'_1, \dots, \mathbf{X}'_{train}]$ and $\tilde{\mathbf{X}} = [\tilde{\mathbf{X}}_1, \dots, \tilde{\mathbf{X}}_{train}]$ in the notation of Sect. 3. We then compute the reconstruction errors of the validation data for each value of $p = 1, \dots, p_{max}$. For the reconstruction, we divide each realization \mathbf{X}_i of the validation data into blocks of length $l \geq p$. A block denotes l states $\mathbf{x}_i^{(j)} = [x_{jl}, \dots, x_{(j+1)l-1}]$, while the next block will be $\mathbf{x}_i^{(j+1)} = [x_{(j+1)l}, \dots, x_{(j+2)l-1}]$. We then compute a reconstruction $\hat{\mathbf{x}}_i^{(j)} = [\hat{x}_{jl}, \dots, \hat{x}_{(j+1)l-1}]$ of this block with the NAR model obtained with SINAR for which we use the last p values of the previous block as starting values. We calculate the relative Euclidean error between reconstruction and data for each block by

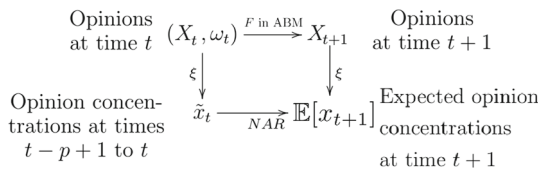
$$err(\hat{\mathbf{x}}_i^{(j)}) = \frac{\|\mathbf{x}_i^{(j)} - \hat{\mathbf{x}}_i^{(j)}\|_F}{\|\mathbf{x}_i^{(j)}\|_F}.$$

Afterwards, we take the mean over all $err(\hat{\mathbf{x}}_i^{(j)})$ to measure the performance of the NAR model.

³ Supposing that there are K basis functions to be used to approximate the space \mathcal{H} , a tensor product basis for the complete space of “delay functions” $\mathbb{Y}^p \rightarrow \mathbb{Y}$ would require K^p functions. Meanwhile, the Mori–Zwanzig formalism does not mix terms from different delays, essentially working on $\bigoplus_{i=1}^p \mathcal{H}$, that is approximated by pK functions.

Since the entries of $\xi(X_t)$ always sum up to 1, information about the percentages of opinions $1, \dots, M - 1$ immediately yields the percentage of opinion M so that we use SINAR to find an NAR model for the evolution of the percentages of the first $M - 1$ opinions only and omit the redundant information $\xi(X)_M$. For the reconstruction error, we compare data about the percentages of only the first $M - 1$ opinions with their reconstructions. This NAR model does not necessarily ensure that the predicted first $M - 1$ percentages stay between 0 and 1 and their sum is at most 1. Since we make short-term predictions only, however, there will at most be only slight deviations from this property.

In the form of the diagram (2.2) from Sect. 2, the Mori–Zwanzig procedure applied to this concept can be described as



Case 1: A Complete Network

For $p_{between} = 1$, the network is complete and there should be no improvement of the prediction by allowing memory terms.

We set $N = 5000$, $T = 300$ and $A_{ij} = 1 \forall i, j$. The number of different opinions is $M = 3$. As coefficients $\alpha_{m'm''}$ we choose

$$\begin{bmatrix} \alpha_{11} & \alpha_{12} & \alpha_{13} \\ \alpha_{21} & \alpha_{22} & \alpha_{23} \\ \alpha_{31} & \alpha_{32} & \alpha_{33} \end{bmatrix} = \begin{bmatrix} 0 & 0.165 & 0.03 \\ 0.03 & 0 & 0.165 \\ 0.165 & 0.03 & 0 \end{bmatrix}$$

As initial percentages we assign values to the $(X_0)_i$ so that $\xi(X_0) = [0.45, 0.1, 0.45]^T$.

As the block length in the validation data, we use $l = 40$. We can already write down the macrodynamics since they are given in (4.1) (see Appendix C.1 for details):

$$\begin{aligned}
 \mathbb{E}[(x_{t+1})_1 | x_t] &= (1 + \alpha_{31} - \alpha_{13})(x_t)_1 + (\alpha_{13} - \alpha_{31})(x_t)_1^2 + (\alpha_{21} - \alpha_{12} - \alpha_{31} + \alpha_{13})(x_t)_1(x_t)_2 \\
 &= 1.135(x_t)_1 - 0.135(x_t)_1^2 - 0.27(x_t)_1(x_t)_2, \\
 \mathbb{E}[(x_{t+1})_2 | x_t] &= (1 + \alpha_{32} - \alpha_{23})(x_t)_2 + (\alpha_{23} - \alpha_{32})(x_t)_2^2 + (\alpha_{12} - \alpha_{21} - \alpha_{32} + \alpha_{23})(x_t)_1(x_t)_2 \\
 &= 0.865(x_t)_2 + 0.135(x_t)_2^2 + 0.27(x_t)_1(x_t)_2.
 \end{aligned}
 \tag{4.3}$$

Inspired by this structure, we choose as basis functions in SINAR

$$[\tilde{\varphi}_1, \dots, \tilde{\varphi}_L](x_t) = [(x_t)_1, (x_t)_2, (x_t)_1^2, (x_t)_2^2, (x_t)_1(x_t)_2]$$

so that

$$\tilde{\Theta}(\tilde{x}_t) = \underbrace{[(x_t)_1, (x_t)_2, (x_t)_1^2, (x_t)_2^2, (x_t)_1(x_t)_2, \dots]}_{\text{Markovian terms as in (4.3)}} \underbrace{[(x_{t-1})_1, (x_{t-1})_2, (x_{t-1})_1^2, (x_{t-1})_2^2, (x_{t-1})_1(x_{t-1})_2, \dots]}_{\text{Memory terms}}^T. \tag{4.4}$$

Since (4.1), resp. (4.3), describe the expected evolution of the percentages and are thus in the form of deterministic models, we omit the noise term ε_{t+1} from (2.9) which we assumed to satisfy $\mathbb{E}[\varepsilon_{t+1}] = 0$.

We create $r = 20$ realizations of which we use 12 for training and the others for validation. We set the sparsity parameter to $\lambda = 0$ and to $\lambda = 0.05$ to test how the accuracy decreases with a sparser model. Since the macrodynamics (4.3) are Markovian, we obtain for the prediction error of the validation data no improvement by allowing memory terms (Fig. 3) for neither the 40- nor the one-step prediction error. Note that the predictions with the sparse NAR model provide slightly better accuracy for large memory depths. This is because small nonzero coefficients for memory terms improve the fit of the training data, but cause errors in the prediction of the validation data, because the macrodynamics are Markovian. Through the sparsity constraint enforced, these nonzero coefficients for memory terms are cut off. The recovered sparse macrodynamics for $p = 1$ reads

$$\begin{aligned} (x_{t+1})_1 &= 1.1353(x_t)_1 - 0.1351(x_t)_1^2 - 0.2709(x_t)_1(x_t)_2, \\ (x_{t+1})_2 &= 0.8655(x_t)_2 + 0.1344(x_t)_2^2 + 0.2699(x_t)_1(x_t)_2, \end{aligned}$$

which is very close to the analytically derived macrodynamics (4.3).

Case 2: A Two-Cluster Network

We now construct a network with $N = 5000$ agents, divided into two clusters of size 2500 each. We set $p_{between} = 0.0001$. Again, $M = 3$ and $\alpha_{m'm''}$ are the same as in case 1. As the starting condition, we let opinions in the first cluster be distributed by $[0.8, 0.1, 0.1]$ and in the second cluster by $[0.1, 0.1, 0.8]$. If the initial percentages in both clusters were equal then the percentages in both clusters would evolve in a quite similar way in parallel so that the macrodynamics would essentially be the same as in the complete network case. With the initial percentages being so different, it is possible that an opinion that is dominant in one cluster at one point in time but only sparsely represented in the other can become popular through the links between agents from different clusters. This will cause the difference in behaviour of the evolution of percentages compared to the complete network.

Morever, in order to derive the Markovian macrodynamics in Eq. (4.1), we needed that the probabilities for an agent i to change its opinion $(X_t)_i$ at time t , which we denoted by $p^t((X_t)_i, m'')$, be independent of i . If the neighbourhoods of different agents are generally different from each other, this is no longer the case. Especially so, if agents are distributed into different clusters, where opinion percentages might

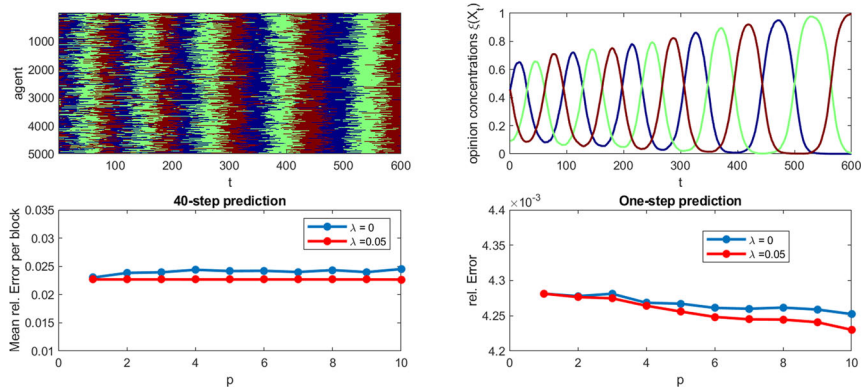


Fig. 3 Results for the complete network. Top left: One realization of the microdynamics. Every column of the graphic represents the opinion of each of the 5000 agents at one point in time. Blue denotes opinion 1, green denotes opinion 2 and red denotes opinion 3. Top right: Corresponding realization of the macrodynamics $\xi(X)$ that represent the percentages of opinions among all agents. We can observe oscillatory behaviour since agents with opinion 1 tend to change their opinion to 2 and analogously from 2 to 3 and from 3 to 1. Bottom: 40-step and one-step relative prediction errors of the NAR models determined by SINAR for different memory depths p with $\lambda = 0$ and $\lambda = 0.05$. As expected, the prediction error does not decrease with higher memory depth than $p = 1$ (Color figure online)

be very different. Thus, we cannot derive Markovian macrodynamics for this case, but in light of the Mori–Zwanzig formalism, we will need memory terms.

To show this, we create $r = 20$ realizations of length $T = 500$ and again use 12 for training, the remaining for validation. As block length, we choose $l = 20$. Memory terms become immediately significant, as the error graphs illustrate (Fig. 4). We use the basis given in (4.4), which has the length $5p$.

The non-sparse and sparse solutions only deviate slightly from each other in their accuracy, but the sparse solution gives a significantly more compact model. For example, for $p = 2$, we obtain for the coefficients $\tilde{\Xi}$

$$\lambda = 0 : \tilde{\Xi} = \begin{bmatrix} 2.04 & 0.03 & -0.07 & -0.08 & 0.02 & -1.05 & -0.02 & 0.07 & 0.07 & -0.02 \\ -0.05 & 1.88 & 0.00 & 0.11 & 0.06 & 0.06 & -0.89 & -0.01 & -0.12 & -0.05 \end{bmatrix}$$

$$\lambda = 0.05 : \tilde{\Xi} = \begin{bmatrix} 1.9691 & 0 & 0 & 0 & -0.9700 & 0 & 0 & 0 & 0 \\ 0 & 1.9662 & 0 & 0 & 0 & -0.9671 & 0 & 0 & 0 \end{bmatrix}$$

so that for $\lambda = 0.05$ the NAR model is given by

$$(x_{t+1})_1 = 1.9691(x_t)_1 - 0.9700(x_{t-1})_1$$

$$(x_{t+1})_2 = 1.9662(x_t)_2 - 0.9671(x_{t-1})_2.$$

For $p = 1$, the NAR model obtained with SINAR ($\lambda = 0.05$) is

$$(x_{t+1})_1 = 1.0094(x_t)_1 - 0.053(x_t)_1(x_t)_2$$

$$(x_{t+1})_2 = 0.9894(x_t)_2 + 0.0574(x_t)_1(x_t)_2.$$

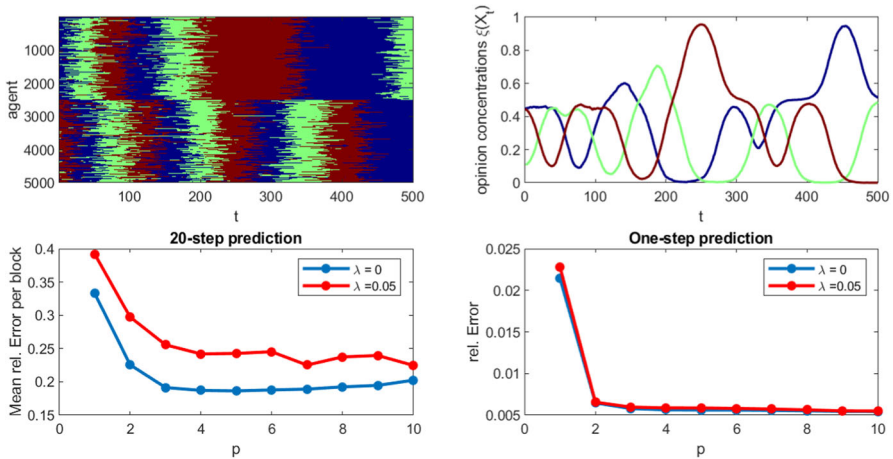


Fig. 4 Results for the two-cluster network. Top left: One realization of the microdynamics. Colours represent opinions as in Fig. 3. Top right: Corresponding realization of the macrodynamics $\xi(X)$. Again there is oscillatory behaviour but also plateaus and short dips as in the red and green graphs at time 25 - 150. This is because at these times one opinion is dominant in one cluster, but not present in the other. Through the links between the clusters, an opinion, that is not present in a cluster but dominant in the other one can be revived, e.g., the blue opinion in the upper cluster. Bottom: 20-step and one-step relative prediction errors of the NAR models determined by SINAR for different memory depths p with $\lambda = 0$ and $\lambda = 0.05$. Memory terms yield a significant decrease in the prediction errors compared to Markovian predictions

With $\lambda = 0$, the obtained NAR model has other terms with nonzero coefficients, but these are small. In Fig. 5, an example for the predictions of opinion percentages in one block using the NAR models with $p = 1, 2$ and 10 is depicted and compared to the corresponding data. As the error graphs in Fig. 4 show already, the predicted percentages come closer to the percentages in the data with increasing memory depth. In order to illustrate why memory terms improve the prediction accuracy, let us imagine for now that there are no links between the clusters. Then, the evolutions of opinion percentages in both clusters run in parallel to each other and are Markovian as derived previously. The opinion percentages in the full network are then given by the averages of the cluster-wise percentages $x_t^{(i)}$, i.e., $x_t = \frac{1}{2}(x_t^{(1)} + x_t^{(2)})$. This means, if we know x_t , then there are various options for what $x_t^{(1)}$ and $x_t^{(2)}$ can be, all of which might result in different values for $x_{t+1}^{(1)}$ and $x_{t+1}^{(2)}$ and thus x_{t+1} . If we are additionally given x_{t-1} , this might yield possible values for $x_{t-1}^{(1)}$ and $x_{t-1}^{(2)}$, which themselves make some of the candidates for $x_t^{(1)}$ and $x_t^{(2)}$ unlikely. Thus, through the information of memory terms we can restrict the options for what the percentages inside each cluster are. We illustrate this in more detail in Appendix C.2.

The links between the clusters have as consequence that within one cluster agents generally do not have identical opinion change probabilities since their neighbourhoods are different. This yields additional need for memory terms since then not even for the macrodynamics in one cluster a Markovian formulation can be derived.

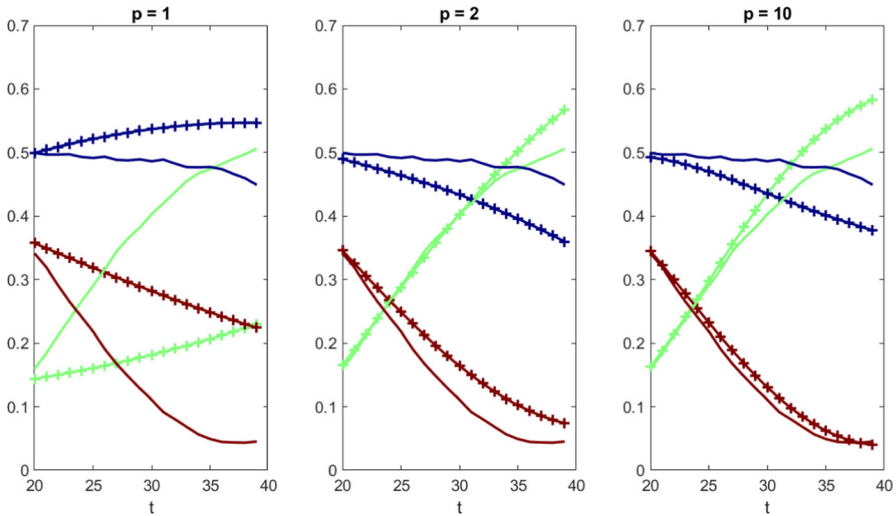


Fig. 5 Opinion percentages over one block of length 20 from the validation data and prediction evolutions with NAR models obtained with SINAR for $p = 1, 2$ and 10 and $\lambda = 0$ (two-cluster network). Percentages from validation data are depicted with thin lines and predicted percentages with lines with crosses. With $p = 1$, the prediction accuracy is poor and improves drastically for $p = 2$. With $p = 10$, the predicted evolutions come even closer to the curves from the validation data

Case 3: A Five-Cluster Network

We repeat the same procedure as with the two-cluster network, but with five clusters of equal size 1000. Again, all agents within a cluster are connected with each other and $p_{between} = 0.0001$. The $\alpha_{m'm''}$ are identical to the ones used in the first two examples. As starting conditions we let opinions in the different clusters be drawn according to different distributions for each cluster. Those distributions are $[0.8, 0.1, 0.1]$, $[0.1, 0.1, 0.8]$, $[0.1, 0.8, 0.1]$, $[0.3, 0.4, 0.3]$ and $[0.5, 0.3, 0.2]$. The evolution of the opinion percentages is now much more irregular compared to the previous examples. The oscillatory behaviour is still present, but the amplitudes differ from time to time. Through the higher number of clusters, more randomness comes into the model since an opinion can be randomly spread from one cluster, where it is dominant, to another one, where it is not dominant, suddenly altering the evolution of percentages in this cluster and thus in the whole network.

We now show that, similar to when we used a two-cluster network, memory terms become important for predictions of the evolution of the microdynamics. This is shown in Fig. 6. Again, the mean relative error per block converges with increasing p . While in the two-cluster network example the performance did not improve visibly with $p > 10$, in this case we can get slightly lower errors for p approaching 20.

For $p = 2$ and $\lambda = 0.05$, we obtain the NAR model

$$\begin{aligned} (x_{t+1})_1 &= 1.8745(x_t)_1 - 0.8748(x_{t-1})_1 \\ (x_{t+1})_2 &= 1.8672(x_t)_2 - 0.8674(x_{t-1})_2. \end{aligned}$$

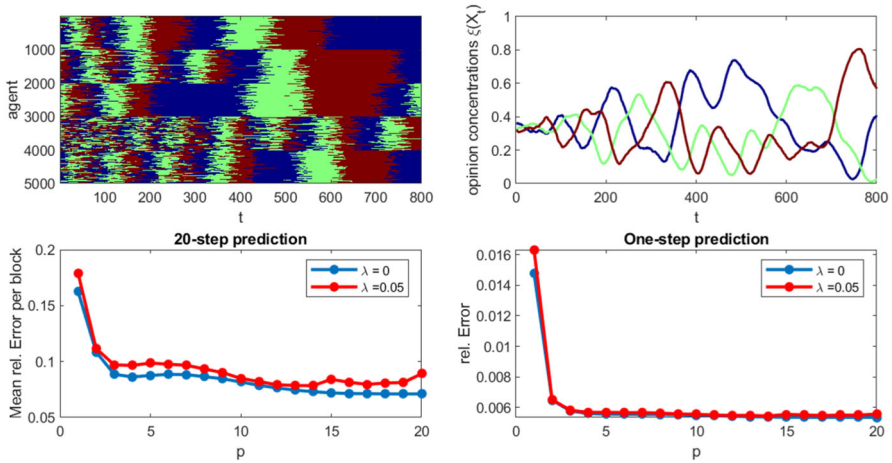


Fig. 6 Results for the five-cluster network. Top left: One realization of the microdynamics. Every column of the graphic represents the opinion of each of the 5000 agents at one point in time. Top right: Corresponding realization of the macrodynamics $\xi(X)$. The behaviour is much more complex than in the first two cases. Bottom: 20-step and one-step relative prediction errors of the NAR models determined by SINAR for different memory depths p with $\lambda = 0$ and $\lambda = 0.05$

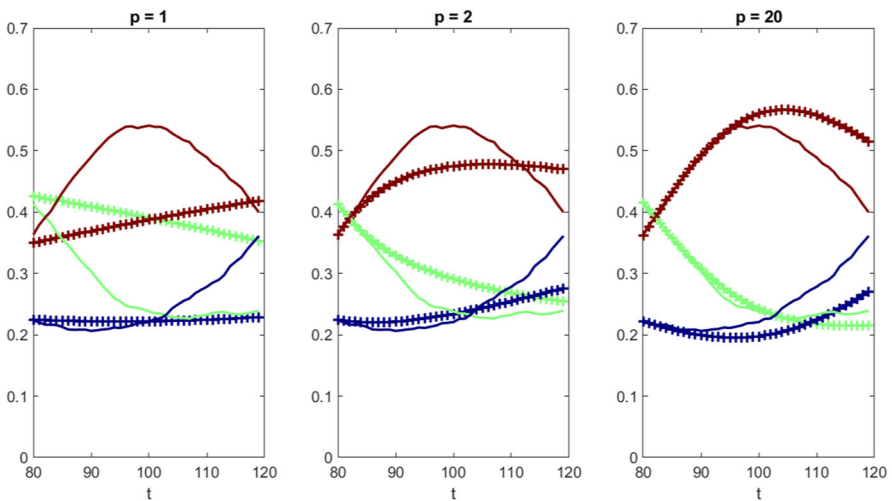


Fig. 7 Opinion percentages over one block of length 40 from the validation data and prediction evolutions with NAR models obtained with SINAR for $p = 1, 2$ and 20 and $\lambda = 0$ (five-cluster network). Percentages from validation data are depicted with thin lines and predicted percentages with lines with crosses. As in the example with a two-cluster network, we can see what the error graphs in Fig. 6 indicate: the predicted evolutions are closer to the validation data with higher memory depths of the NAR model

For $p > 2$, the models show increasing complexity, e.g., for $p = 3$:

$$\begin{aligned}(x_{t+1})_1 &= 1.4662(x_t)_1 - 0.1188(x_t)_2 + 0.0552(x_t)_1^2 + 0.1318(x_t)_1(x_t)_2 \\ &\quad + 0.2309(x_{t-1})_2 \\ &\quad - 0.1899(x_{t-1})_1(x_{t-1})_2 - 0.2021(x_{t-1})_2^2 - 0.4658(x_{t-2})_1 \\ &\quad - 0.1060(x_{t-2})_2 \\ &\quad + 0.1206(x_{t-2})_1^2 + 0.0644(x_{t-2})_2^2(x_{t+1})_2 \\ &= 1.3157(x_t)_2 - 0.3161(x_{t-2})_2.\end{aligned}$$

Again, we show as an example the predictions of percentages for one block of length 40 with memory depths 1, 2 and 10 (Fig. 7). As in the example with the two-cluster network, we can see that a higher memory depth indeed increases the prediction accuracy for the evolution of the opinion percentages in the short term, i.e., for predictions of length 20 resp. 40. Plus, enforcing the sparsity constraint with the parameter λ in SINAR set to 0.05 yields significantly sparser models, while the prediction accuracy only suffered slightly.

5 Discussion

In this article, we have summarized how the evolution of observations of a dynamical system can be derived through the Mori–Zwanzig formalism and how this can result in a nonlinear autoregressive model with memory. For the determination of model parameters, we have used methodology from data-driven system identification methods, inspired by SINDy (Brunton et al. 2016a). We could then extend SINDy to SINAR which identifies sparse nonlinear autoregressive (NAR) models from data, thus deploying a common system identification method for non-Markovian systems.

We applied this to an agent-based model (ABM) that simulates the dynamics of opinion changes in a population. Assuming that all agents are equally strongly influenced by all other agents in the population, we showed that for the prediction of the percentages of opinions within the population memory terms are not necessary. However, for incomplete networks, this is no longer the case. Our methodology enabled us to make more accurate predictions for the percentages of opinions among the agents when the population of agents was defined by clusters with little influence between them. Additionally, sparse models obtained from enforcing a sparsity constraint in the estimation of NAR models in SINAR gave almost equally good prediction accuracy as the non-sparse ones, while yielding far simpler models. In the context of opinion dynamics, such sparse models permit to point out more clearly which opinions impact which others and how.

The following challenges have yet to be addressed:

- In our methodology, we have assumed a noise term resulting from Mori–Zwanzig that was zero mean. This allowed us to omit it when making predictions of the expected value of the opinion percentages. This simplifying assumption does not need to be true, and one could try to derive a more accurate representation for the

noise term. As a result of this simplifying assumption, the NAR models we considered were deterministic, even for non-deterministic microdynamics. Introduction of explicit noise in the NAR models, e.g., by extending the approach outlined in Klus et al. (2020), could improve their (statistical) predictive capacities.

- One could additionally choose a different projection P in the Mori–Zwanzig formalism. The choice of an orthogonal projection on a finite set of basis functions explicitly yielded an NAR model. The right projection for a given system could inspire an optimal choice of basis functions, e.g., such that the memory depth is minimal.
- We have derived models that are stationary, i.e., do not change over time. Since the assumption of an equilibrium distribution over states of the microdynamics might not always hold, coefficients of the NAR model may become time-dependent. One could use a regime switching model as in Horenko (2011) that fixes coefficients for a time interval before changing them to other fixed values when the macrodynamics show certain behaviour, e.g., coefficients might be different depending on which opinion is dominating.

A MATLAB toolbox for the experiments done in Sect. B and Appendix 4 is provided under <https://github.com/nwulkow/OpinionDynamicsModelling>.

Acknowledgements PK and CS acknowledge support by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s Excellence Strategy—The Berlin Mathematics Research Center MATH+ (EXC-2046/1, project ID: 390685689). NW thanks Luzie Helfmann, Jan-Hendrik Niemann and Alexander Sikorski for helpful discussions on the subject of opinion dynamics.

Funding Open Access funding enabled and organized by Projekt DEAL.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

A Technical Details on the Mori–Zwanzig Equation and SINAR

A.1 The Derivation of the Mori–Zwanzig Equation

We show here how to derive Eq. (2.5) from the Dyson formula in Sect. 2.2. The Dyson formula states

$$\mathcal{K}^{t+1} = \sum_{k=0}^t \mathcal{K}^{t-k} P\mathcal{K}(Q\mathcal{K})^k + (Q\mathcal{K})^{k+1}.$$

Application of both sides of the equation to ξ and evaluation at the initial value X_0 yield

$$\begin{aligned}
 (\mathcal{K}^{t+1}\xi)(X_0) &= \sum_{k=0}^t \mathcal{K}^{t-k}[P\mathcal{K}(Q\mathcal{K})^k\xi](X_0) + (Q\mathcal{K})^{t+1}\xi(X_0) \\
 \text{which results in: } \xi(X_{t+1}) &= \sum_{k=0}^t [P\mathcal{K}(Q\mathcal{K})^k\xi](X_{t-k}) + (Q\mathcal{K})^{t+1}\xi(X_0) \\
 \text{Setting } \rho^k := (Q\mathcal{K})^k\xi \text{ yields: } \xi(X_{t+1}) &= \sum_{k=0}^t [P\mathcal{K}\rho^k](X_{t-k}) + \rho^{t+1}(X_0) \\
 &= \sum_{k=0}^t [P(\rho^k \circ F)](X_{t-k}) + \rho^{t+1}(X_0).
 \end{aligned}
 \tag{A.1}$$

We can replace X_{t-k} by x_{t-k} in the last step because the application of P to a function makes this function depend only on the relevant variables. We explicitly used the parentheses around the operator $P\mathcal{K}\rho^k$ and its equivalent formulations to indicate that P is a projection operator that works on the function $\mathcal{K}\rho^k$.

Since $\rho^0 = \xi$, we obtain that $P(\rho^0 \circ F) = P(\xi \circ F)$. This is usually referred to as the *optimal prediction* term since it is the best Markovian approximation of $\xi(X_{t+1})$, i.e., the best approximation of $\xi(X_{t+1})$ that only uses $\xi(X_t)$. The sum in the last row of (A.1) starting at $k = 1$ is referred to as the *memory terms*, since these terms use information from previous values of $\xi(X)$. The term $\rho^{t+1}(X_0)$ depending on the full state X_0 and not on the projection $\xi(X_0)$, is often called *noise*, because one does not have explicit access to it and can often only treat it as a stochastic influence.⁴ In total, the last row of (A.1) is called the *Mori–Zwanzig equation*.

Substituting the definition of P as the orthogonal projection onto basis functions as in (2.3), we obtain

$$\begin{aligned}
 P(\rho^k \circ F)(x_{t-k}) &= \varphi(x_{t-k})\langle\varphi, \varphi\rangle^{-1}\langle\varphi, \rho^k \circ F\rangle \\
 &= \underbrace{\varphi(x_{t-k})}_{\in \mathbb{R}^{m \times L}} \underbrace{\langle\varphi, \varphi\rangle^{-1}}_{\in \mathbb{R}^{L \times L}} \underbrace{\int_{\mathbb{X}} \underbrace{\varphi(\xi(X))^T}_{\in \mathbb{R}^{L \times m}} \underbrace{\rho^k(F(X))}_{\in \mathbb{Y} \subset \mathbb{R}^m} d\mu(X)}_{\in \mathbb{R}^L} \\
 &=: \varphi(x_{t-k})h_k \in \mathbb{R}^m
 \end{aligned}
 \tag{A.2}$$

with vector-valued coefficients $h_k = \langle\varphi, \varphi\rangle^{-1} \int_{\mathbb{X}} \varphi(\xi(X))^T \rho^k(F(X))d\mu(X)$.

A.2 Translations Between the Model Forms (2.8) and (2.9)

We show here how to translate a model in the form of (2.8) into the form of (2.9) and vice versa. Starting in the form of (2.8), we suppose we have chosen basis functions

⁴ It accumulates unobserved effects as witnessed by the complement projector Q . Note that it is expected to decay fast, if the system mixes strongly (in the sense that \mathcal{K} has a small spectral radius on the set of functions perpendicular to the constant function, which in turn is assumed to lie in the range of P). In this sense, the term “noise” refers to negligible correlation to variables x_{t-k} that contribute strongly to $\xi(X_{t+1})$.

$\varphi = [\varphi_1, \dots, \varphi_L] \in \mathbb{R}^{m \times L}$ and $h_k = [(h_k)_1, \dots, (h_k)_L] \in \mathbb{R}^L$. This gives

$$\varphi(x_{t-k})h_k = \sum_{i=1}^L (h_k)_i \varphi_i(x_{t-k}).$$

Let us choose $\tilde{\varphi} = [\varphi_1^T, \dots, \varphi_L^T]^T \in \mathbb{R}^{mL}$, set $H_k^{(i)} = h_i I_{m \times m} \in \mathbb{R}^{m \times m}$ and define $H_k = [H_k^{(1)}, \dots, H_k^{(L)}] \in \mathbb{R}^{m \times mL}$. Then

$$\begin{aligned} H_k \tilde{\varphi}(x_{t-k}) &= [H_k^{(1)}, \dots, H_k^{(L)}] \begin{bmatrix} \varphi_1(x_{t-k}) \\ \vdots \\ \varphi_L(x_{t-k}) \end{bmatrix} = \sum_{i=1}^L H_k^{(i)} \varphi_i(x_{t-k}) \\ &= \sum_{i=1}^L (h_k)_i \varphi_i(x_{t-k}) \\ &= \varphi(x_{t-k})h_k. \end{aligned}$$

Thus, we can express (2.8) in the form of (2.9) by imposing the restriction on the matrices H_k that they have the form $H_k = h_k I_{m \times m}$. Note that we have simply modified the forms in which the dynamics are expressed, but not generated a different model structure.

For the backward direction, suppose we have chosen scalar-valued basis functions $\tilde{\varphi}_1, \dots, \tilde{\varphi}_K$ and determined matrix-valued coefficients $H_k \in \mathbb{R}^{m \times K}$. Then we can bring (2.9) into the form of (2.8) by setting $L = mK$, defining φ as the Kronecker product $\varphi = I_{m \times m} \otimes [\tilde{\varphi}_1, \dots, \tilde{\varphi}_K]$, i.e.,

$$\varphi(x) = \begin{bmatrix} \tilde{\varphi}_1(x) & \dots & \tilde{\varphi}_K(x) & 0 & \dots & \dots & 0 \\ 0 & \dots & 0 & \tilde{\varphi}_1(x) & \dots & \tilde{\varphi}_K(x) & 0 & \dots & 0 \\ \vdots & & & & & & & & \\ 0 & & \dots & & & & \tilde{\varphi}_1(x) & \dots & \tilde{\varphi}_K(x) \end{bmatrix} \in \mathbb{R}^{m \times mK},$$

and using mK -dimensional coefficients

$$h_k = [(H_k)_{11}, \dots, (H_k)_{1K}, \dots, (H_k)_{m1}, \dots, (H_k)_{mK}]^T.$$

Then,

$$\begin{aligned} \varphi(x_{t-k})h_k &= \varphi(x_{t-k}) \begin{bmatrix} (H_k)_{11} \\ \vdots \\ (H_k)_{mK} \end{bmatrix} = \begin{bmatrix} (H_k)_{11} & \dots & (H_k)_{1K} \\ \vdots & & \vdots \\ (H_k)_{m1} & \dots & (H_k)_{mK} \end{bmatrix} \begin{bmatrix} \tilde{\varphi}_1(x_{t-k}) \\ \vdots \\ \tilde{\varphi}_K(x_{t-k}) \end{bmatrix} \\ &= H_k \tilde{\varphi}(x_{t-k}). \end{aligned}$$

A.3 Relation Between SINDy, SINAR, DMD and AR

The diagram in Fig. 1 sketches how system identification methods from different contexts are related. With DMD, SINDy, SINAR and AR models in mind, one can observe that in all of them, a minimization problem of the same form is solved: given are data matrices \mathbf{X} and \mathbf{X}' which contain data points of the realization of a (possibly memory-exhibiting) dynamical system that are shifted from each other by one time step. Then one tries to find a connection between both through a transformation of \mathbf{X} which is multiplied with a coefficient matrix by solving (omitting possible sparsity constraints)

$$\Xi = \arg \min_{\Xi} \|\mathbf{X}' - \Xi \Theta(\mathbf{X})\|_F$$

In DMD, one tries to find a linear and Markovian connection between x_t and x_{t+1} , i.e., $\Theta(x) = x$. In SINDy, \mathbf{X} is transformed in a possibly nonlinear way in order to explain the evolution of systems for which a linear model might be inaccurate.

Linear AR models look for a linear connection between a fixed number of past values of the system and its next value. The columns of \mathbf{X} , in this case, contain not just data points of the system but sequences of data points of a fixed length. More precisely,

In DMD, one minimizes
$$\sum_{t=0}^{T-1} \|x_{t+1} - \Xi x_t\|_2$$

In AR models, define $\tilde{x}_t = [x_t^T, \dots, x_{t-p+1}^T]^T$ and minimize
$$\sum_{t=p-1}^{T-1} \|x_{t+1} - \tilde{\Xi} \tilde{x}_t\|_2$$

Since in DMD one maps time-shifted versions of the same coordinates onto each other (i.e., x_t to x_{t+1}), let us augment the AR minimization problem to $\sum_{t=p-1}^{T-1} \|\tilde{x}_{t+1} - C \tilde{x}_t\|_2$. Then $\tilde{\Xi}$ is equal to the upper m rows of C , while the lower $m(p-1)$ rows of C have simple structure copying the associated rows from \tilde{x}_t (C is a so-called *companion matrix*). In this way, the AR problem is equivalent to the DMD problem with states from the Hankel matrix defined in (3.5). In Arbabi and Mezic (2017), the authors discuss Hankel-DMD to extract properties of the Koopman operator of a system from observational data. In doing so, they essentially fit an AR model.

In the same fashion, SINAR is the delay-embedded counterpart to SINDy and brings together SINDy and AR models in the sense that one seeks a possibly nonlinear connection between past values of the system and subsequent ones.

A.4 Definition of the Conditional Expectation

Let states $X \in \mathbb{X}$ be distributed according to μ . Let us define for $\xi \in \mathcal{G}$ the level sets $L_x := \{X \in \mathbb{X} : \xi(X) = x\}$. Then, through the coarea formula (Federer 1996), the

expectation of a function $g \in \mathcal{G}$ with $g \in L^1(\mathbb{X})$ can be written as

$$\begin{aligned} \mathbb{E}_\mu[g(X)] &= \int_{\mathbb{X}} f(X) d\mu(X) \\ &= \int_{\xi(\mathbb{X})} \int_{L_x} f(X) \mu(X) \det(\nabla \xi(X)^T \nabla \xi(X))^{-\frac{1}{2}} dx d\sigma_x(X) \end{aligned}$$

where σ_x is the Hausdorff measure on L_x . Then, the conditional expectation of $f(X)$ given that $\xi(X) = x$ is (see, e.g., Bittracher et al. (2018))

$$\mathbb{E}_\mu[g(X) \mid \xi(X) = x] = \frac{1}{\Gamma(x)} \int_{L_x} g(X) \mu(X) \det(\nabla \xi(X)^T \nabla \xi(X))^{-\frac{1}{2}} d\sigma_x(X),$$

where $\Gamma(x)$ is a normalization constant.

A.5 Determination of Coefficients of Linear AR Models

A linear autoregressive model with zero-mean Gaussian noise has the form

$$x_{t+1} = \sum_{i=0}^{p-1} H_i x_{t-i} + \varepsilon_{t+1}, \quad \varepsilon_{t+1} \sim \mathcal{N}(0, \Sigma^T \Sigma).$$

The best linear unbiased estimator (BLUE) (Plackett 1949; Baksalary and Kala 1981) for the H_i is the least squares minimizer $\tilde{\Xi} = [H_0, \dots, H_{p-1}]$, given by

$$\tilde{\Xi} = \arg \min_{\tilde{\Xi}=[H_0, \dots, H_{p-1}]} \|\mathbf{X}' - \Xi \tilde{\mathbf{X}}\|_F,$$

where $\tilde{\mathbf{X}}$ and \mathbf{X}' are defined as in (3.5).

Omitting the sparsity constraint, SINAR solves the problem

$$\tilde{\Xi} = \arg \min_{\tilde{\Xi}} \|\mathbf{X}' - \tilde{\Xi} \Theta(\tilde{\mathbf{X}})\|_F.$$

If $\Theta(\tilde{x}) = \tilde{x}$, then this is precisely the least squares method for linear autoregressive models.

A.6 Covariance of Noise Terms of NAR Models

Assuming a relation of the form

$$x'_t = \Xi \Theta(x_t) + \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(0, \Sigma^T \Sigma),$$

we find that

$$Cov(x'_t - \Xi \Theta(x_t)) = Cov(\varepsilon_t).$$

An unbiased estimator for the covariance of a random variable y is the statistical covariance

$$\bar{\Sigma} = \frac{1}{T-1} \sum_{t=1}^T (y_t - \bar{y})(y_t - \bar{y})^T$$

where $\bar{y} = \frac{1}{T} \sum_{t=1}^T y_t$.

In order to estimate the covariance matrix of noise terms ε_{t+1} in Eq. (2.9), one has to substitute x'_t by x_{t+1} and $\Xi\Theta(x_t)$ by $\sum_{k=0}^{p-1} H_k \tilde{\varphi}(x_{t-k})$ to derive the form of Eq. (2.9).

Subsequently y has to be substituted by $x_{t+1} - \sum_{k=0}^{p-1} H_k \tilde{\varphi}(x_{t-k})$ and we can calculate the statistical covariance of ε_{t+1} in (2.9).

B Example: Application of SINAR to an Extended Hénon System

We demonstrate here the emergence of memory terms in the case of inaccessible variables in the sense of the Mori–Zwanzig formalism by means of an example of a dynamical system and use SINAR to detect an NAR model that reconstructs the dynamics.

B.1 The Classical Hénon System and an Extension

The classical Hénon system (Hénon 1976) describes a two-dimensional system that is one of the most famous examples for systems with chaotic behaviour, i.e., where slightly deviated initial conditions lead to a significantly different trajectory. The dynamical system is given by

$$\begin{aligned} x_{t+1} &= 1 - ax_t^2 + y_t \\ y_{t+1} &= bx_t, \end{aligned}$$

where a, b are fixed parameters. As we can observe, y_t is nothing more than a scaled and time-delayed version of x_t . We now consider x as the relevant and y as the irrelevant variable; this means in the Mori–Zwanzig formalism the space \mathcal{H} is given by all functions depending on only x . We can then still express the evolution of x exactly with dependence on the past two values of x by plugging in the equation for y_{t+1} into the equation for x_{t+1} :

$$x_{t+1} = 1 - ax_t^2 + bx_{t-1}.$$

Let us now consider an extended version of the Hénon system

$$\begin{aligned} x_{t+1} &= 1 - ax_t^2 + y_t \\ y_{t+1} &= bx_t + cy_t \end{aligned} \tag{B.1}$$

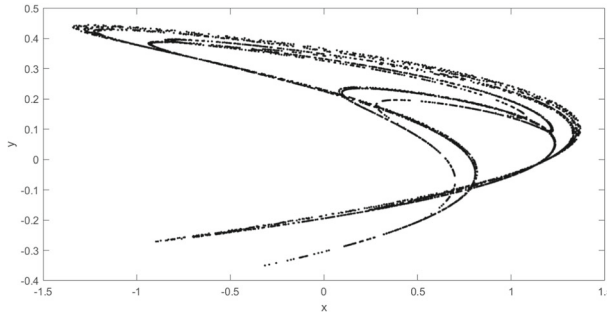


Fig. 8 Trajectory of length 5000 of the two-dimensional extended Hénon system (B.1) with $a = 1.3, b = 0.3, c = 0.3$ and initial values $x_0 = y_0 = 0$. The first 1000 states are omitted here so that the trajectory has time to converge towards the attractor

whose dynamical behaviour is visualized in Fig. 8. Now y is more than only a scaled and time-delayed version of x . If we try to express x_t only in dependence of its own past terms and without values of y , then we do not get a system with a finite memory depth, but with an infinite one:

$$\begin{aligned}
 x_{t+1} &= 1 - ax_t^2 + bx_{t-1} + cy_{t-1} \\
 &= 1 - ax_t^2 + bx_{t-1} + cbx_{t-2} + c^2y_{t-2} \\
 &= 1 - ax_t^2 + bx_{t-1} + cbx_{t-2} + c^2b_{t-3} + c^3y_{t-3} \\
 &= 1 - ax_t^2 + \sum_{j=1}^t c^{j-1}bx_{t-j} + c^{t+1}y_0,
 \end{aligned} \tag{B.2}$$

which can be quickly shown by induction on t .

We have hereby derived an equation of the form of the Mori–Zwanzig equation (2.5) for this simple example: the term $1 - ax_t^2$ is the optimal prediction, i.e., the Markovian approximation using the relevant variables x_t . The sum

$$\sum_{j=1}^t c^{j-1}bx_{t-j}$$

contains the memory terms depending on past values of x and the term $c^t y_0$ is the noise term with information about the irrelevant, or for us inaccessible, variable y .

B.2 Reconstructing the Extended Hénon System with SINAR

We now apply the SINAR algorithm to data originating from a trajectory of the extended Hénon system and demonstrate the increase in performance by using memory terms compared to applying the usual Markovian SINDy algorithm.

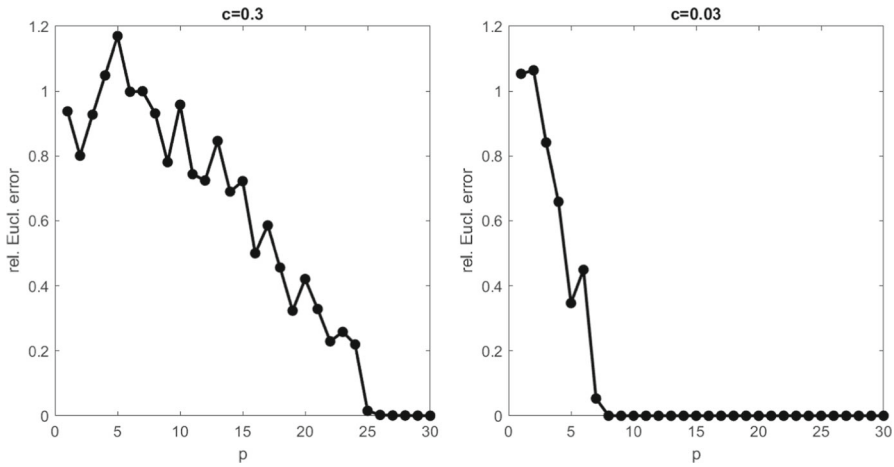


Fig. 9 Relative error of validation $err(\hat{X}')$ for SINAR on visible variable x of the extended Hénon system for two different values of c with different memory depths p on the x axis. The prediction accuracy improves with increasing memory depth. Results based on SINAR with $\lambda = 0$. As parameters in the extended Hénon system, we chose $a = 1.3, b = 0.3$ and $c = 0.3$ (left) resp. $c = 0.03$ (right). For every value of p , the same 80 time steps were taken into account for the reconstruction error

We set as parameters $a = 1.3, b = 0.3, c = 0.3$ and initial values $x_0 = y_0 = 0$. Then, for example, the exact model up to a memory depth of 3 in Eq. (B.2) is

$$x_{t+1} = 1 - 1.3x_t^2 + 0.3x_{t-1} + 0.09x_{t-2} + 0.027x_{t-3} + \mathcal{O}(c^3).$$

As basis functions, we choose monomials of the time-delayed coordinates up to second order without mixed terms between different delays,

$$\tilde{\Theta}(\tilde{x}_t) = \begin{bmatrix} 1 \\ x_t^2 \\ x_t \\ \vdots \\ x_{t-p+1} \end{bmatrix}.$$

Short-Term Predictions

We now generate a trajectory of length $T = 2000$ out of which we erase the first 1000 steps to give the trajectory time to converge to the attractor. We then use the first T_{train} data points for training and the remaining $1000 - T_{train}$ for validation. With the training data, we determine coefficients $\tilde{\Xi}$ for the basis functions in $\tilde{\Theta}$ with SINAR for different memory depths p and compute reconstructions $\hat{x}_{T_{train}+1}, \dots, \hat{x}_{1000}$ of $x_{T_{train}+1}, \dots, x_{1000}$ using Eq. (3.7) with initial values $x_{T_{train}-p+1}, \dots, x_{T_{train}}$. In essence, we recover the coefficients of the forms a resp. $c^{j-1}b$ from Eq. (B.2) until

$j = p - 1$ and recompute values of the extended Hénon system with the recovered coefficients. As error measure we use the relative Euclidean prediction error

$$\text{err}(\hat{\mathbf{X}}') = \frac{\|\mathbf{X}' - \hat{\mathbf{X}}'\|_F}{\|\mathbf{X}'\|_F} \quad (\text{B.3})$$

where $\mathbf{X}' = [x_{T_{\text{train}}+1}, \dots, x_{1000}]$ denotes data points from the original trajectory and $\hat{\mathbf{X}}' = [\hat{x}_{T_{\text{train}}+1}, \dots, \hat{x}_{1000}]$ data points from the reconstructed trajectory.

Although all coefficients are recovered up to an error of smaller than 10^{-14} when we use 800 time steps for training, the reconstruction becomes inaccurate after around 100 time steps which underlines the strongly chaotic nature of the system, i.e., small deviations at one point in time causing significant deviations in the long-term behaviour. We thus use 920 time steps for training and only 80 time steps for validation to investigate how the relative Euclidean reconstruction error depends on the memory depth. Below we discuss how the attractor of the system is recovered using much longer reconstructions.

We see in Fig. 9 how the relative Euclidean prediction error decreases for increasing memory depth p . Predicted was the evolution of x with data about x . It is interesting to note how large a memory depth is necessary to get an accurate prediction for x when $c = 0.3$ (Fig. 9 (left)). The chaotic nature of the system yields that even coefficients of the form bc^j for $j = 27$ have to be taken into account. Of course, for smaller c such as $c = 0.03$, memory terms in (B.2) decay quicker and a memory depth of $p = 8$ is sufficient to yield an accurate prediction as shown in Fig. 9 (right). For the full system (x, y) , the system is Markovian and the prediction error is unsurprisingly very small even for $p = 1$.

Attractor Reconstruction

Although large deviations between original and reconstructed trajectories of x_t occur after around 100 time steps, both trajectories remain on roughly the same set of points. We quantify this by the Hausdorff distance between the two-dimensional delay embeddings (see definition in Appendix B.3) of the original trajectory and each reconstructed trajectory. The Hausdorff distance denotes the maximal minimal distance of members of one set of points to another set. In other words, the Hausdorff distance between two sets is 0 if the sets are equal and big if there is a point in one set which is far away from all points in the other set.

We make predictions of 3000 time steps based on coefficients that were obtained with SINAR on data of 1000 time steps. In Fig. 11 are depicted the two-dimensional delay embeddings of the original trajectory of x and the reconstructed trajectories for $p = 1, 2, 5, 10$ and $p = 30$. There we see how already for $p = 2$ the original and reconstructed attractors look much more similar compared to $p = 1$. Figure 10 shows the Hausdorff distances for different memory depths. Similar to the relative Euclidean prediction error, the distance decreases with increasing p . The remaining error is

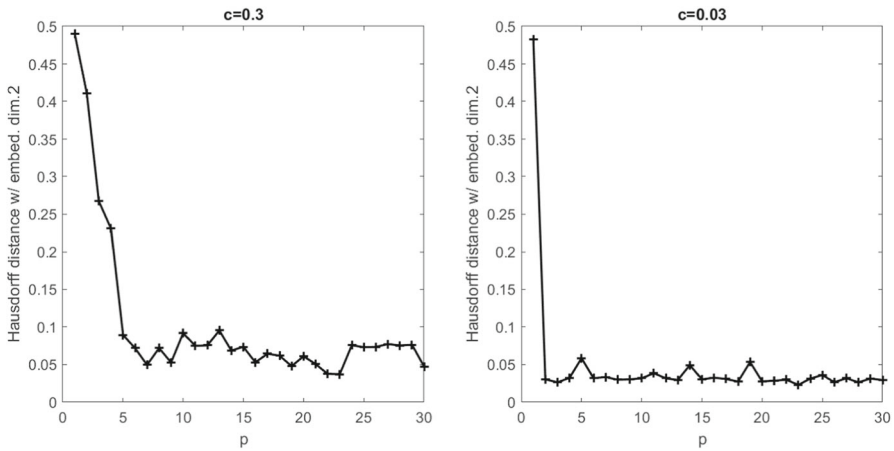


Fig. 10 Hausdorff distances between original and reconstructed attractors with 3000 points of two-dimensional delay embeddings of x for two different values of c with different memory depths p on the x axis. Results based on SINAR with $\lambda = 0$ with parameters in the extended Hénon system $a = 1.3, b = 0.3$ and $c = 0.3$ (left) resp. $c = 0.03$ (right)

due to the fact that the complicated geometry of the attractor is hard to approximate uniformly well with a finite set of points (Fig. 11).⁵

B.3 Hausdorff Distance of Delay Embedding of Trajectories

The Hausdorff distance between two non-empty compact sets measures the maximal minimal distance a point from one set has to the other set. It is commonly used to compare attractors of dynamical systems. The lower the Hausdorff distance between two sets, the more similar they are. From two trajectories $\mathbf{X}' = [x_0, \dots, x_T]$ and $\hat{\mathbf{X}}' = [\hat{x}_0, \dots, \hat{x}_T]$, we construct the delay embeddings with embedding depth p as

$$\mathcal{D}_p(\mathbf{X}') = \left[\begin{bmatrix} x_{p-1} \\ \vdots \\ x_0 \end{bmatrix}, \begin{bmatrix} x_p \\ \vdots \\ x_1 \end{bmatrix}, \dots \right], \quad \mathcal{D}_p(\hat{\mathbf{X}}') = \left[\begin{bmatrix} \hat{x}_{p-1} \\ \vdots \\ \hat{x}_0 \end{bmatrix}, \begin{bmatrix} \hat{x}_p \\ \vdots \\ \hat{x}_1 \end{bmatrix}, \dots \right].$$

We then calculate their Hausdorff distance as

$$\max \left(\max_{x \in \mathcal{D}_p(\mathbf{X}')} \min_{\hat{x} \in \mathcal{D}_p(\hat{\mathbf{X}}')} \|x - \hat{x}\|_2, \max_{\hat{x} \in \mathcal{D}_p(\hat{\mathbf{X}}')} \min_{x \in \mathcal{D}_p(\mathbf{X}')} \|x - \hat{x}\|_2 \right).$$

⁵ Coverage of a two-dimensional object of diameter 2 by 3000 points results in a mesh size $\approx 2/\sqrt{3000} \approx 0.03$. This is the same order of magnitude as the error we observe.

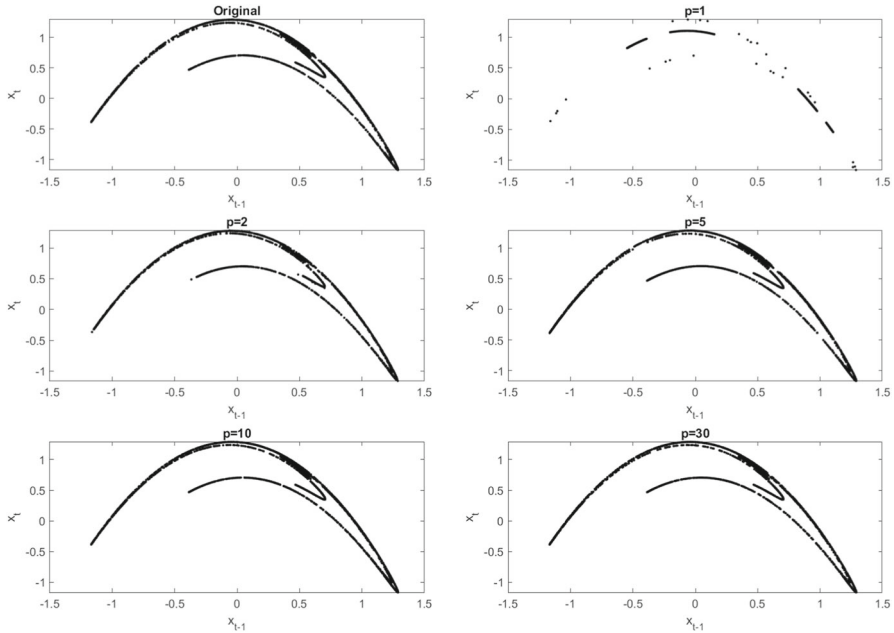


Fig. 11 Two-dimensional delay embedded attractors with 3000 points of the extended Hénon system with $a = 1.3, b = 0.3, c = 0.3$. Original (upper left) and reconstructed ones based SINAR with $\lambda = 0$. For $p \geq 2$, differences are difficult to see, but exist as the Hausdorff distances in Fig. 10 indicate

C Details on Expected Opinion Dynamics

C.1 Derivation of Eq. (4.3)

With $m = 3$ opinions, Eq. (4.1) reads

$$\begin{aligned} (x_{t+1})_1 &= (x_t)_1 + (\alpha_{21} - \alpha_{12})(x_t)_1(x_t)_2 + (\alpha_{31} - \alpha_{13})(x_t)_1(x_t)_3 \\ (x_{t+1})_2 &= (x_t)_2 + (\alpha_{12} - \alpha_{21})(x_t)_1(x_t)_2 + (\alpha_{32} - \alpha_{23})(x_t)_2(x_t)_3 \\ (x_{t+1})_3 &= (x_t)_3 + (\alpha_{13} - \alpha_{31})(x_t)_1(x_t)_3 + (\alpha_{23} - \alpha_{32})(x_t)_2(x_t)_3. \end{aligned}$$

Using $(x_t)_3 = 1 - (x_t)_1 - (x_t)_2$, we get

$$\begin{aligned} (x_{t+1})_1 &= (x_t)_1 + (\alpha_{21} - \alpha_{12})(x_t)_1(x_t)_2 + (\alpha_{31} - \alpha_{13})(x_t)_1(1 - (x_t)_1 - (x_t)_2) \\ (x_{t+1})_2 &= (x_t)_2 + (\alpha_{12} - \alpha_{21})(x_t)_1(x_t)_2 + (\alpha_{32} - \alpha_{23})(x_t)_2(1 - (x_t)_1 - (x_t)_2). \end{aligned}$$

Rearranging gives

$$\begin{aligned} (x_{t+1})_1 &= (1 + \alpha_{31} - \alpha_{13})(x_t)_1 + (\alpha_{13} - \alpha_{31})(x_t)_1^2 \\ &\quad + (\alpha_{21} - \alpha_{12} - \alpha_{31} + \alpha_{13})(x_t)_1(x_t)_2 \\ (x_{t+1})_2 &= (1 + \alpha_{32} - \alpha_{23})(x_t)_2 + (\alpha_{23} - \alpha_{32})(x_t)_2^2 \\ &\quad + (\alpha_{12} - \alpha_{21} - \alpha_{32} + \alpha_{23})(x_t)_1(x_t)_2. \end{aligned}$$

This is Eq. (4.3).

C.2 Representations of Uncoupled Expected Two-Cluster Dynamics

In this subsection, we discuss the derivation of NAR models for a network which consists of two equally sized clusters without links between them. Having derived the expected dynamics for a complete network in Eq. (4.1), we assume for now that the expected dynamics are identical with the true dynamics in order to investigate the macrodynamics if the agents behave perfectly as expected. We then get Markovian deterministic dynamics that describe the evolution of opinion percentages in each cluster. Their means are the opinion percentages in the whole network. The derivation of an NAR model for this property is analytically challenging but numerical results suggest certain structures of the macrodynamics dependent on the initial percentages.

Macrodynamics inside the clusters.

Since the clusters represent complete networks of their own, we obtain for the opinion percentages $x_t^{(i)}$ inside each cluster

$$\begin{aligned} (x_{t+1}^{(i)})_1 &= (1 + \alpha_{31} - \alpha_{13})(x_t^{(i)})_1 + (\alpha_{13} - \alpha_{31})(x_t^{(i)})_1^2 + (\alpha_{21} - \alpha_{12} - \alpha_{31} + \alpha_{13})(x_t^{(i)})_1(x_t^{(i)})_2 \\ (x_{t+1}^{(i)})_2 &= (1 + \alpha_{32} - \alpha_{23})(x_t^{(i)})_2 + (\alpha_{23} - \alpha_{32})(x_t^{(i)})_2^2 + (\alpha_{12} - \alpha_{21} - \alpha_{32} + \alpha_{23})(x_t^{(i)})_1(x_t^{(i)})_2. \end{aligned} \tag{C.1}$$

With $x_t = \frac{1}{2}(x_t^{(1)} + x_t^{(2)})$ and denoting $a = \alpha_{31} - \alpha_{13}$, $b = \alpha_{21} - \alpha_{12} - \alpha_{31} + \alpha_{13}$, $c = \alpha_{32} - \alpha_{23}$, $d = \alpha_{12} - \alpha_{21} - \alpha_{32} + \alpha_{23}$, this gives

$$\begin{aligned} (x_{t+1})_1 &= (1 + a) \underbrace{\frac{1}{2}((x_t^{(1)})_1 + (x_t^{(2)})_1)}_{(x_t)_1} - \frac{a}{2}((x_t^{(1)})_1^2 + (x_t^{(2)})_1^2) + \frac{b}{2}((x_t^{(1)})_1(x_t^{(1)})_2 + (x_t^{(2)})_1(x_t^{(2)})_2) \\ (x_{t+1})_2 &= (1 + c) \underbrace{\frac{1}{2}((x_t^{(1)})_2 + (x_t^{(2)})_2)}_{(x_t)_2} - \frac{c}{2}((x_t^{(1)})_2^2 + (x_t^{(2)})_2^2) + \frac{d}{2}((x_t^{(1)})_1(x_t^{(1)})_2 + (x_t^{(2)})_1(x_t^{(2)})_2). \end{aligned}$$

Even making the simplifying assumption that $a = -c$ and $b = -d = -2a$ as is the case for the coefficients we chose for the examples, we arrive at

$$\begin{aligned} (x_{t+1})_1 &= (1 + a) \underbrace{\frac{1}{2}((x_t^{(1)})_1 + (x_t^{(2)})_1)}_{(x_t)_1} - \frac{a}{2}((x_t^{(1)})_1^2 + (x_t^{(2)})_1^2) - a((x_t^{(1)})_1(x_t^{(1)})_2 + (x_t^{(2)})_1(x_t^{(2)})_2) \\ (x_{t+1})_2 &= (1 - a) \underbrace{\frac{1}{2}((x_t^{(1)})_2 + (x_t^{(2)})_2)}_{(x_t)_2} + \frac{a}{2}((x_t^{(1)})_2^2 + (x_t^{(2)})_2^2) + a((x_t^{(1)})_1(x_t^{(1)})_2 + (x_t^{(2)})_1(x_t^{(2)})_2). \end{aligned}$$

From this, it seems impossible to find a closed Markovian expression for x_t . In order to understand why memory terms should help to express the evolution of x_t , note the following: given x_{t-1} and x_t , we could now find $x_{t-1}^{(1)}$ and $x_{t-1}^{(2)}$ so that these equations would yield those values for $x_t^{(1)}$ and $x_t^{(2)}$ whose average is x_t . This set of pairs of $x_t^{(1)}$ and $x_t^{(2)}$ would significantly be limited compared to all pairs which have this x_t as their

average. From these $x_t^{(i)}$, we could compute subsequent values $x_{t+1}^{(i)}$. Hence, we would have gained a more precise estimate of $x_t^{(1)}$ and $x_t^{(2)}$ and thus of x_{t+1} . In the stochastic ABM, the evolution of x_t is originally stochastic if it represents the percentages of opinions of agents. Hence, one would not search for the $x_{t-1}^{(i)}$ that exactly yield x_t , but rather make this argument in terms of probabilities. We would then get different probabilities for the $x_t^{(i)}$ dependent on what x_{t-1} is.

Simplified example: Linear dynamics inside the clusters.

Of course, a closed expression for the evolution of x_{t+1} that depends only on memory terms of x_t and not on the $x_t^{(i)}$ is desirable. However, the analytical derivation of such an expression seems out of reach. Thus, as an example for much simpler macrodynamics inside each cluster, we illustrate how one can find a closed expression for the mean of two linear dynamics. For this, let

$$\begin{aligned} x_{t+1}^{(1)} &= \lambda_1 x_t^{(1)} \\ x_{t+1}^{(2)} &= \lambda_2 x_t^{(2)} \end{aligned}$$

and

$$x_t = \frac{1}{2}(x_t^{(1)} + x_t^{(2)}).$$

Thus,

$$\begin{aligned} x_t^{(i)} &= \lambda_i^t x_0^{(i)}, \quad i = 1, 2 \\ \text{and } x_t &= \frac{1}{2}(\lambda_1^t x_0^{(1)} + \lambda_2^t x_0^{(2)}). \end{aligned}$$

Then one can observe that

$$x_{t+1} = \frac{(\lambda_1 + \lambda_2)}{2} x_t - \frac{\lambda_1 \lambda_2}{2} x_{t-1}$$

since

$$\begin{aligned} &\frac{(\lambda_1 + \lambda_2)}{2} x_t - \frac{\lambda_1 \lambda_2}{2} x_{t-1} \\ &= \frac{1}{2}[(\lambda_1 + \lambda_2)(\lambda_1^t x_0^{(1)} + \lambda_2^t x_0^{(2)}) - \lambda_1 \lambda_2 (\lambda_1^{t-1} x_0^{(1)} + \lambda_2^{t-1} x_0^{(2)})] \\ &= \frac{1}{2}[(\lambda_1^{t+1} x_0^{(1)} + \lambda_2^{t+1} x_0^{(2)}) + \lambda_1 \lambda_2^t x_0^{(2)} + \lambda_2 \lambda_1^t x_0^{(1)} - \lambda_2 \lambda_1^t x_0^{(1)} - \lambda_1 \lambda_2^t x_0^{(2)}] \\ &= \frac{1}{2}(\lambda_1^{t+1} x_0^{(1)} + \lambda_2^{t+1} x_0^{(2)}) = \frac{1}{2}(x_{t+1}^{(1)} + x_{t+1}^{(2)}) = x_{t+1}. \end{aligned}$$

Numerical results with symmetric initial percentages.

For the macrodynamics (C.1) of opinion percentages in a two-cluster network, we have not derived such a closed expression analytically. However, we can see numerically that almost exact models can be derived for a memory depth of $p = 2$ if we impose *symmetric* starting conditions, i.e., initial percentages that fulfill

$$(x_0^{(1)})_1 = (x_0^{(2)})_2 = 1 - 2(x_0^{(1)})_2, \quad (x_0^{(2)})_1 = (x_0^{(1)})_2 = 1 - 2(x_0^{(2)})_2.$$

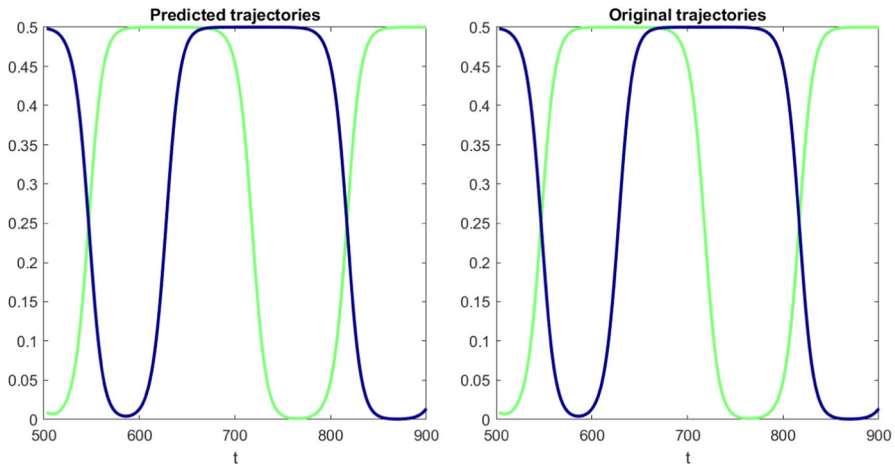


Fig. 12 Original trajectories for initial percentages in (C.2) and predicted trajectories with the NAR model (C.3)

To illustrate this, we create trajectories of length $T = 900$ of the deterministic dynamics (C.1) with initial percentages

$$(x_0^{(1)})_1 = 0.8, \quad (x_0^{(1)})_2 = 0.1, \quad (x_0^{(2)})_1 = 0.1, \quad (x_0^{(2)})_2 = 0.8 \quad (C.2)$$

and $a = 0.135$ which is also the case in the examples in Sect. 4.

From the first 500 time steps of the resulting $x_t = \frac{1}{2}(x_t^{(1)} + x_t^{(2)})$, we estimate the NAR model (with $\lambda = 0$ in SINAR)

$$\begin{aligned} (x_{t+1})_1 &= 1.21(x_t)_1 - 0.65(x_t)_2 + 0.27(x_t)_1^2 + 0.54(x_t)_1(x_t)_2 \\ &\quad - 0.26(x_{t-1})_1 + 0.71(x_{t-1})_2 - 0.17(x_{t-1})_1^2 - 0.10(x_{t-1})_2^2 \\ &\quad - 0.54(x_{t-1})_1(x_{t-1})_2 \\ (x_{t+1})_2 &= -0.82(x_t)_1 - 1.31(x_t)_2 - 0.27(x_t)_2^2 - 0.54(x_t)_1(x_t)_2 \\ &\quad + 0.68(x_{t-1})_1 - 0.16(x_{t-1})_2 - 0.30(x_{t-1})_1^2 - 0.03(x_{t-1})_2^2 \\ &\quad + 0.54(x_{t-1})_1(x_{t-1})_2. \end{aligned} \quad (C.3)$$

With this model, we reconstruct the remaining 400 time steps in the data by computing a trajectory of length 400 with starting values given by x_{499} and x_{500} (Fig. 12). The relative Euclidean error between both trajectories amounts to $2.4 \cdot 10^{-7}$. For the one-step prediction, i.e., mapping every two values x_{t-1} and x_t to x_{t+1} with the above model, the error is $1.5 \cdot 10^{-14}$. For larger memory depths, there is no improvement in prediction accuracy. This suggests that for these specific initial conditions the macrodynamics can be reproduced with memory depth $p = 2$.

Numerical results with non-symmetric initial percentages.

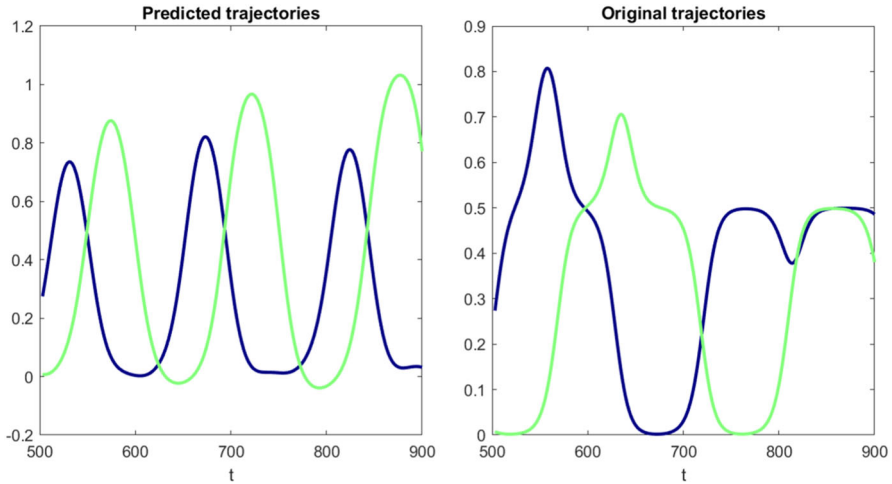


Fig. 13 Original trajectories for initial percentages in (C.4) and predicted trajectories with the NAR model (C.2)

For other initial percentages, we get quite different coefficients that significantly decrease the influence of the second-order terms $(x_t)_1^2$, $(x_t)_2^2$ and $(x_t)_1(x_t)_2$. Let

$$(x_0^{(1)})_1 = 0.7, \quad (x_0^{(1)})_2 = 0.2, \quad (x_0^{(2)})_1 = 0.1, \quad (x_0^{(2)})_2 = 0.8. \quad (C.4)$$

Then for $p = 2$, in the same manner ($\lambda = 0$), we obtain the model

$$\begin{aligned} (x_{t+1})_1 &= 2.09(x_t)_1 - 0.01(x_t)_2 - 0.09(x_t)_1^2 + 0.01(x_t)_2^2 - 0.15(x_t)_1(x_t)_2 \\ &\quad - 1.09(x_{t-1})_1 + 0.01(x_{t-1})_2 + 0.09(x_{t-1})_1^2 - 0.02(x_{t-1})_2 \\ &\quad + 0.15(x_{t-1})_1(x_{t-1})_2 \\ (x_{t+1})_2 &= -0.04(x_t)_1 - 1.90(x_t)_2 - 0.01(x_t)_1^2 - 0.08(x_t)_2^2 + 0.15(x_t)_1(x_t)_2 \\ &\quad + 0.04(x_{t-1})_1 - 0.90(x_{t-1})_2 - 0.08(x_{t-1})_2^2 - 0.16(x_{t-1})_1(x_{t-1})_2. \end{aligned}$$

The original trajectories and the trajectories obtained from this model are depicted in Fig. 13.

The one-step prediction error improves for memory depths larger than $p = 2$ (Fig. 14). Since with NAR models obtained from the trajectories for these initial percentages, the predicted trajectories diverge, the full prediction error is not shown.

In summary, for a network that consists of two clusters which are uncoupled but fully connected internally, the expected macrodynamics are given by the mean of the expected intra-cluster dynamics. Assuming the dynamics to have no variance and hence to be deterministic, given in Eq. (C.1), with symmetric initial percentages, a memory depth of 2 is enough for us to generate an almost exact NAR model for the macrodynamics. However, for non-symmetric initial percentages, the ensuing best-fitting NAR models with the basis functions we use are not accurate in the long term.

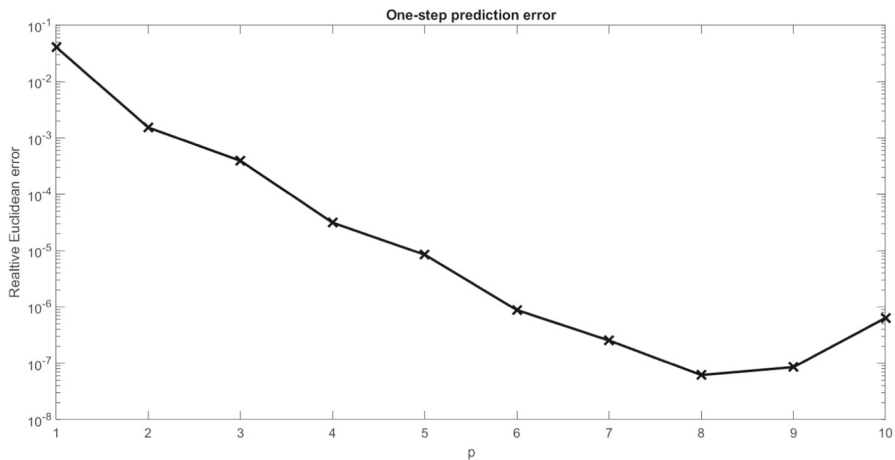


Fig. 14 One-step prediction error for the NAR models obtained from trajectories of x_t with initial percentages of the $x_t^{(i)}$ as given in (C.4)

This seems to be in part due to the fact that for non-symmetric initial percentages, the trajectories show more complex behaviour which no longer consists of periodic oscillations, but is rather more irregular. This could cause the best-fitting NAR models to then be dominated by linear terms. Results about to which degree one can analytically derive NAR models for both symmetric and non-symmetric initial percentages require further research.

References

- Aho, K., Derryberry, D., Peterson, T.: Model selection for ecologists: the worldviews of AIC and BIC. *Ecology* **95**(3), 631–636 (2014)
- An, H., Huang, F.: The geometrical ergodicity of nonlinear autoregressive models. *Stat. Sin.* **6**, 943–956 (1996)
- Anderson, B.D.O., Ye, M.: Recent advances in the modelling and analysis of opinion dynamics on influence networks. *Int. J. Autom. Comput.* **16**, 129–149 (2019)
- Arbabi, H., Mezic, I.: Ergodic theory, dynamic mode decomposition and computation of spectral properties of the koopman operator. *SIAM J. Appl. Dyn. Syst.* **4**(16), 2096–2126 (2017)
- Baksalary, J., Kala, R.: Simple Least Squares estimation versus best linear unbiased prediction. *J. Stat. Plan. Inference* **2**(5), 147–151 (1981)
- Banisch, S.: From microscopic heterogeneity to macroscopic complexity in the contrarian voter model. *Adv. Complex Syst.* **12**, 1450025 (2014)
- Banisch, S.: *Markov Chain Aggregation for Agent-Based Models*, vol. 1. Springer, Berlin (2016)
- Banisch, S., Lima, R., Araújo, T.: Agent based models and opinion dynamics as Markov chains. *Soc. Netw.* **34**, 549–561 (2011)
- Berryman, A.A.: The origins and evolution of predator-prey theory. *Ecology* **73**(5), 1530–1535 (1992)
- Billings, S.: *Nonlinear System Identification: NARMAX Methods in the Time, Frequency, and Spatio-Temporal Domains*, vol. 1. Wiley, Hoboken (2013)
- Bittracher, A., Koltai, P., Klus, S., Banisch, R., Dellnitz, M., Schütte, C.: Transition manifolds of complex metastable systems. *J. Nonlinear Sci.* **28**, 471–512 (2018)
- Böhme, G.A., Gross, T.: Fragmentation transitions in multistate voter models. *Phys. Rev. E* **85**, 066117 (2012)

- Bolzern, P., Colaneri, P., Nicolao, G.: Opinion influence and evolution in social networks: a markovian agents model. *Automatica* **100**, 11 (2017)
- Boschia, G., Cammarotaa, C., Kühna, R.: Opinion dynamics with memory: how a society is shaped by its own past. [arXiv:1909.12590](https://arxiv.org/abs/1909.12590), 09 (2019)
- Bowman, G.R., Pande, V.S., Noé, F.: *An Introduction to Markov State Models and Their Application to Long Timescale Molecular Simulation*, vol. 1. Springer, Berlin (2014)
- Brockwell, P.J., Davis, R.A.: *Time Series: Theory and Methods*, vol. 2. Springer, Berlin (1991)
- Brunton, S., Brunton, B., Proctor, J., Kaiser, E., Kutz, J.: Chaos as an intermittently forced linear system. *Nat. Commun.* **8**, 1–9 (2016)
- Brunton, S.L., Proctor, J.P.L., Kutz, J.N.: Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proc. Natl. Acad. Sci.* **113**(15), 3932–3937 (2016)
- Brunton, S.L., Proctor, J.P.L., Kutz, J.N.: Sparse Identification of Nonlinear Dynamics with Control (SINDyC). *IFAC-PapersOnLine Issue* **18**(49), 710–715 (2016)
- Castellano, C., Fortunato, S., Loreto, V.: Statistical physics of social dynamics. *Rev. Mod. Phys.* **81**, 591–646 (2009)
- Chen, G., Duan, X., Friedkin, N., Bullo, F.: Social power dynamics over switching and stochastic influence networks. *IEEE Trans. Autom. Control* **04**, 1 (2018)
- Chorin, A.J., Hald, O.H., Kupferman, R.: Optimal prediction and the Mori–Zwanzig representation of irreversible processes. *Proc. Natl. Acad. Sci.* **97**(7), 2968–2973 (2000)
- Chorin, A.J., Hald, O.H., Kupferman, R.: Optimal prediction with memory. *Phys. D* **166**, 239–257 (2002)
- Davis, R., Zang, P., Zheng, T.: Sparse vector autoregressive modeling. *J. Comput. Graph. Stat.* **30**, 1077–1096 (2012)
- De, A., Bhattacharya, S., Bhattacharya, P., Ganguly, N., Chakrabarti, S.: Learning linear influence models in social networks from transient opinion dynamics. *ACM Trans. Web* **13**, 1–33 (2019)
- Devroye, L., Györfi, L., Lugosi, G.: *A probabilistic theory of pattern recognition*, vol. 31. Springer, Berlin (2013)
- Federer, H.: *Geometric Measure Theory*, vol. 1. Springer, Berlin (1996)
- Fujita, A., Sato, J., Garay, M., Yamaguchi, R., Miyano, S., Sogayar, M., Ferreira, C.: Modeling gene expression regulatory networks with the sparse vector autoregressive model. *BMC Syst. Biol.* **1**, 39 (2007)
- Gilani, F., Giannakis, D., Harlim, J.: Kernel-based prediction of non-markovian time series, 07 (2020)
- Hansen, P., D. O’leary, : The use of the L-curve in the regularization of discrete ill-posed problems. *SIAM J. Sci. Comput.* **14**, 1487–1503 (1993)
- Hénon, M.: A two-dimensional mapping with a strange attractor. *Commun. Math. Phys.* **50**(1), 69–77 (1976)
- Hijón, C., Español, P., Vanden-Eijnden, E., Delgado-Buscalioni, R.: Mori–Zwanzig formalism as a practical computational tool. *Faraday Discuss.* **144**, 301–22 (2010). discussion 323
- Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural Comput.* **9**, 1735–80 (1997)
- Horenko, I.: On analysis of nonstationary categorical data time series: dynamical dimension reduction, model selection, and applications to computational sociology. *Multiscale Model. Simul.* **9**, 1700–1726 (2011)
- Horenko, I., Hartmann, C., Schütte, C., Noe, F.: Data-based parameter estimation of generalized multidimensional Langevin processes. *Phys. Rev. E* **76**, 016706 (2007)
- Jedrzejewski, A., Sznajd-Weron, K.: Impact of memory on opinion dynamics. *Phys. A Stat. Mech. Appl.* **505**, 03 (2018)
- Jennings, N., Sycara, K., Wooldridge, M.: A roadmap of agent research and development. *Auton. Agents Multi-Agent Syst.* **1**, 7–38 (1998)
- Jovanovic, M., Schmid, P., Nichols, J.: Sparsity-promoting dynamic mode decomposition. *Phys. Fluids* **26**, 1–22 (2013)
- Kaiser, E., Kutz, J.N., Brunton, S.L.: Sparse identification of nonlinear dynamics for model predictive control in the low-data limit. *Proc. R. Soc. A* **474**, 20180335 (2018)
- Klimek, P., Lambiotte, R., Thurner, S.: Opinion formation in laggard societies. *Europhys. Lett. EPL* **82**, 1–5 (2007)
- Klus, S., Nüske, F., Peitz, S., Niemann, J.-H., Clementi, C., Schuette, C.: Data-driven approximation of the Koopman generator: model reduction, system identification, and control. *Phys. D* **406**, 132416 (2020)
- Kondrashov, D., Chekroun, M.D., Ghil, M.: Data-driven non-Markovian closure models. *Phys. D Nonlinear Phenom.* **297**, 33–55 (2015)
- Konishi, S., Kitagawa, G.: *Information Criteria and Statistical Modeling*, vol. 1. Springer, Berlin (2008)

- Koopman, B.O.: Hamiltonian systems and transformation in Hilbert space. *Proc. Natl. Acad. Sci.* **17**(5), 315–318 (1931)
- Laubenbacher, R., Jarrach, A. S., Mortveit, H. S., Ravi, S.: *Agent Based Modeling, Mathematical Formalism for*, pp. 160–176. Springer New York, New York, NY, (2009)
- Lei, H., Baker, N.A., Li, X.: Data-driven parameterization of the generalized Langevin equation. *Proc. Natl. Acad. Sci.* **113**(50), 14183–14188 (2016)
- Li, X., Chu, W.: The Mori-Zwanzig formalism for the derivation of a fluctuating heat conduction model from molecular dynamics. *Commun. Math. Sci.* **17**, 539–563 (2017)
- Li, Q., Braunstein, L., Wang, H., Shao, J., Stanley, H., Havline, S.: Non-consensus opinion models on complex networks. *J. Stat. Phys.* **151**, 10 (2012)
- Lin, K.K., Lu, F.: Data-driven model reduction, Wiener projections, and the Mori–Zwanzig formalism. [arXiv:1908.07725v1](https://arxiv.org/abs/1908.07725v1), (2019)
- Lu, F., Maggioni, M., Tang, S., Zhong, M.: Nonparametric inference of interaction laws in systems of agents from trajectory data. *Proc. Natl. Acad. Sci.* **116**, 06 (2019)
- Misra, A.K.: A simple mathematical model for the spread of two political parties. *Nonlinear Anal. Model. Control*, 2012, No. 3 **17**, 343–354 (2012)
- Moussaïd, M., Kämmer, J., Analytis, P., Neth, H.: Social influence and the collective dynamics of opinion formation. *PLoS One* **8**, e78433 (2013)
- Nardini, C., Kozma, B., Barrat, A.: Who’s talking first? Consensus or lack thereof in coevolving opinion formation models. *Phys. Rev. Lett.* **100**, 158701 (2008)
- Pan, S., Duraisamy, K.: Long-time predictive modeling of nonlinear dynamical systems using neural networks. *Complexity* **1–26**, 2018 (2018)
- Plackett, R.: A Historical Note on the Method of Least Squares. *Biometrika*, No. 3/4 **36**, 458–460 (1949)
- Raftery, A.E.: A model for high-order Markov chains. *J. R. Stat. Soc. Ser. B (Methodological)* No. 3 **47**, 528–539 (1985)
- Ravazzi, C., Hojjatnia, S., Lagoa, C., Dabbene, F.: Randomized opinion dynamics over networks: Influence estimation from partial observations. In: *Proceedings of the IEEE Conference on Decision and Control*, pp. 2452–2457. Institute of Electrical and Electronics Engineers Inc., January (2019)
- Schmid, P., Sesterhenn, J.: Dynamic mode decomposition of numerical and experimental data. *J. Fluid Mech.* **656**, 11 (2008)
- Sîrbu, A., Loreto, V., Servodio, V., Tria, F.: *Opinion Dynamics: Models, Extensions and External Effects*, pp. 363–401. 05 (2017)
- Stangor, C.: *Social Groups in Action and Interaction*, vol. 2. Routledge, London (2015)
- Sugihara, G., May, R.M.: Nonlinear forecasting as a way of distinguishing chaos from measurement error in time series. *Nature* **344**, 734–741 (1990)
- Takens, F.: Detecting strange attractors in turbulence. **898**, 366–381 (1981)
- Tibshirani, R.: Regression shrinkage and selection via the Lasso. *J. R. Stat. Soc. Ser. B (Methodological)*, no. 1 **58**, 267–288 (1996)
- Tu, J.H., Rowley, C.W., Luchtenburg, D.M., Brunton, S.L., Kutz, J.N.: On dynamic mode decomposition: theory and applications. *J. Comput. Dyn.* **1**, 391–421 (2014)
- Tuyen, L.: A higher order Markov model for time series forecasting. *Int. J. Appl. Math. Stat.* **57**, 1–18 (2018)
- Venkataramani, S.C., Venkataramani, R.C., Restrepo, J.M.: Dimension reduction for systems with slow relaxation. *J. Stat. Phys.* **167**, 892–933 (2017)
- Williams, M., Kevrekidis, I., Rowley, C.: A data-driven approximation of the Koopman operator: extending dynamic mode decomposition. *J. Nonlinear Sci.* **25**, 1307–1346 (2014)
- Wu, X., Wai, H.-T., Scaglione, A.: Estimating social opinion dynamics models from voting records. *IEEE Trans. Signal Process.* **04**, 1 (2018)
- Xia, H., Wang, H., Xuan, Z.: Opinion dynamics: a multidisciplinary review and perspective on future research. *IJKSS* **2**, 72–91 (2011)
- Zhu, Y., Dominy, J., Venturi, D.: On the estimation of the Mori–Zwanzig memory integral. *J. Math. Phys.* **59**, 103501 (2018)
- Zwanzig, R.: *Nonequilibrium Statistical Mechanics*. Oxford University Press, Oxford (2001)