

# LINEARIZING VISUAL PROCESSES WITH DEEP GENERATIVE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

This work studies the problem of modeling non-linear visual processes by leveraging deep generative architectures for learning linear, Gaussian models of observed sequences. We propose a joint learning framework, combining a *multivariate autoregressive model* and *deep convolutional generative networks*. After justification of theoretical assumptions of linearization, we propose an architecture that allows Variational Autoencoders and Generative Adversarial Networks to simultaneously learn the non-linear observation as well as the linear state-transition model from a sequence of observed frames. Finally, we demonstrate our approach on conceptual toy examples and dynamic textures.

## 1 INTRODUCTION

While classification of image and video with Convolutional Neural Networks (CNN) is becoming an established practice, unsupervised learning and generative modeling remain to be challenging problems in deep learning. A generative model of a visual process enables the possibility of generating sequences of video frames such that the appearance as well as the dynamics approximately resemble the original training process without copying it. This procedure is typically referred to as video generation (Vondrick & Torralba (2017); De Souza et al. (2017)) or video synthesis (Liu et al. (2017b)). More technically, this means that in addition to a suitable probability model for the individual frames, a probabilistic description for the frame-to-frame transition is also necessary. Analysis and reproduction of visual processes simplifies considerably, if this transition can be described as a *multivariate autoregressive* (MAR) model, i.e., as a combination of linear transformations and Gaussian noise. For instance, linear transformations are easily invertible and by means of spectral analysis, it can be studied how such a process behaves in the long term.

Realistically, most frame transitions in real-world visual processes unlikely are linear functions. Nevertheless, unsupervised learning has come up with many approaches to fit MAR models to real-world processes, for instance by using linear low-rank approximations, as proposed by Doretto et al. (2003), or sparse approximations of the frames, as proposed by Wei et al. (2017), or applying the kernel trick to them (Chan & Vasconcelos (2007)).

The success of Generative Adversarial Networks (GAN) introduced by Goodfellow et al. (2014) and Variational Autoencoders (VAE) introduced by Kingma & Welling (2013) has led to an increased interest in deep generative learning and it seems natural to apply such techniques to sequential processes. We approach this idea from the perspective of *linearization*, in order to keep the model as simple as possible. In an analogous way as physicists transforming non-linear differential equations into linear ones by means of an appropriate change of variables, our approach is to *learn* latent representations of visual processes, such that the latent state-to-state transition can be described by an MAR model. To this end, we *jointly* learn a non-linear observation and a linear state transition function by introducing a *dynamic layer* that can be used in conjunction with deep generative architectures such as GANs and VAEs.

## NOTATION

In this paper, letters that appear both in roman and italicized type refer to random variables and realizations thereof, respectively. If not stated otherwise, expected values of more than one random variables assume statistical independence.

## 2 RELATED WORK

The deep predictive coding network (Chalasanani & Principe (2013)), can be considered one of the earliest approaches of combining MAR models with deep learning. More recently, such a combination was studied by Liu et al. (2018). The work by Goroshin et al. (2015) deals with linearizing transformations under uncertainty via neural networks. It resembles our work in that we also focus on *representation learning* rather than on a particular type of process. However, these works do not employ VAEs or GANs.

By contrast, Watter et al. (2015) combine Linear Dynamic Systems (LDSs) with VAEs. However, the focus of their work is on control rather than on synthesis. Furthermore, their model is *locally* linear and the transition distribution is modeled in the variational bound, whereas we model it as a separate layer. This also is the main difference to the work of Krishnan et al. (2015), in which VAEs are used as Kalman Filters. The work of Johnson et al. (2016) combines VAEs with linear dynamic models for forecasting images in video sequences. Mathematically, it is a well-thought out approach, but proposes a training objective that is considerably more complex than the one proposed in this work. GANs in combination with MAR models have been studied in experiments of Han et al. (2017), but the MAR model was learned separately from the GAN.

Theoretical groundwork regarding learned visual transformations has been done by Cohen & Welling (2014; 2016); Hinton et al. (2011) and Memisevic (2013). In particular, the *dynamic layer* proposed in this work bears resemblance to techniques employed in image-to-image translation (Liu et al. (2017a); Zhu et al. (2017)).

The synthesis of video dynamics by means of neural networks has been discussed, among others by Vondrick et al. (2016); Xue et al. (2016); Srivastava et al. (2015) and Xie et al. (2017). We would like to emphasize the difference between video synthesis which is discussed in our work and the *prediction* of video frames as studied, for instance, by Mathieu et al. (2015). While the former refers to the problem of finding a *probabilistic* model for the spatial and temporal behavior of a visual process, the latter refers to finding a *deterministic* mapping from a set of previous frames to one or several future frames to come. As a consequence, a video synthesis model needs to take care of long-term behavior. Additionally, the probabilistic nature of video synthesis makes it considerably harder to evaluate the generated frames, since a unique ground truth cannot be provided, ruling out classical quantitative quality measures such as mean squared error or structured similarity.

Finally, the core contribution of this work is a combination of neural networks with Markov processes. This has been the subject of many works in the recent past. For a broad overview of results in this field, the reader is referred to Chapter 20 of the book *Deep Learning* by Goodfellow et al. (2016).

## 3 VISUAL PROCESSES

### 3.1 DYNAMIC SYSTEMS

The *Dynamic Texture* (DT) model by Doretto et al. (2003) has popularized LDSs in the modeling of visual processes. Typically, an LDS is of the following form

$$\begin{aligned} \mathbf{h}_{t+1} &= \mathbf{A}\mathbf{h}_t + \mathbf{v}_t, \\ \mathbf{y}_t &= \bar{\mathbf{y}} + \mathbf{C}\mathbf{h}_t, \end{aligned} \tag{1}$$

where  $\mathbf{h}_t \in \mathbb{R}^n$  is the low dimensional *state space variable* at time  $t$ ,  $\mathbf{A} \in \mathbb{R}^{n \times n}$  the *state transition matrix*,  $\mathbf{y}_t \in \mathbb{R}^d$  the observation at time  $t$  and  $\mathbf{C} \in \mathbb{R}^{d \times n}$  the *observation matrix*. The vector  $\bar{\mathbf{y}} \in \mathbb{R}^d$  represents a constant offset in the observation space. The input term  $\mathbf{v}_t$  is modeled as i.i.d. zero-mean Gaussian noise, and is independent of  $\mathbf{h}_t$ .

Real-world visual processes are often highly non-linear and non-Gaussian, hence are considerably harder to work with. In the following, we define a more general dynamic model that is applicable to a broad class of visual processes. To keep the problem tractable, we make the assumption of two properties that were already implicitly assumed in the classical DT model.

**Assumption 1.** *The visual processes of interest are stationary first-order Markov processes.*

The stationarity property guarantees that synthesis of new sequences will not diverge and the Markov property facilitates the statistical inference. However, we believe that the method presented in this work can be generalized to Markov processes of higher orders. We refer to Appendix A.2 for details.

Consider the following generic dynamic equation of observations  $\mathbf{y}$  in an  $n$ -dimensional manifold  $\mathbb{M} \subset \mathbb{R}^d$ .

$$\mathbf{y}_{t+1} = \psi_{\mathbf{v}_t}(\varphi(\mathbf{y}_t)). \quad (2)$$

Again,  $\mathbf{y}_t \in \mathbb{R}^d$  denotes the observation at time  $t$ . The self map  $\varphi : \mathbb{M} \rightarrow \mathbb{M}$  models the predicted frame transition and describes the deterministic part of the dynamics. In this work, it is assumed to be differentiable. Furthermore, the function  $\psi_{\mathbf{v}_t} : \mathbb{M} \rightarrow \mathbb{M}$  describes a displacement by  $\mathbf{v}_t$  in the tangent space  $\mathcal{T}_{\mathbf{y}_t}\mathbb{M}$  of  $\mathbb{M}$  at  $\mathbf{y}_t$ , followed by a retraction onto  $\mathbb{M}$ . The displacement is sampled from i.i.d. zero-mean Gaussian noise. Eq. (1) is a special case of Eq. (2), where  $\varphi(\mathbf{y}) = \mathbf{C}\mathbf{A}\mathbf{C}^+(\mathbf{y} - \bar{\mathbf{y}}) + \bar{\mathbf{y}}$  and  $\psi_{\mathbf{v}}(\mathbf{y}) = \mathbf{y} + \mathbf{C}\mathbf{v}$ .

The model Eq. (2) describes a much broader class of visual processes than Eq. (1). However, unlike model Eq. (2), the state transition in the linear model Eq. (1) enables straightforward prediction, generation, and analysis of observations. It is thus of great interest to find a model that linearizes real-world visual processes, so that in some latent state space representation, the state transition admits the MAR model of the first line of Eq. (1). Specifically, this work focuses on the following non-linear dynamic system model, i.e., a linear state transition and a non-linear observation mapping

$$\begin{aligned} \mathbf{h}_{t+1} &= \mathbf{A}\mathbf{h}_t + \mathbf{v}_t, \\ \mathbf{y}_t &= \mathbf{C}(\mathbf{h}_t), \end{aligned} \quad (3)$$

where  $n < d$  and  $\mathbf{C} : \mathbb{R}^n \rightarrow \mathbb{M} \subset \mathbb{R}^d$  is differentiable. Transforming a non-linear model to this form makes video synthesis as easy as sampling from autoregressive noise.

It is worth noticing that the model in Eq. (3) is not unique with respect to changes of basis in the state space (Doretto et al. (2003); Afsari & Vidal (2013)). However, if  $\mathbf{C}$  is implemented via a neural network, we can ensure that it accounts for a possible change of basis, e.g. by adding a linear layer to the input of the network. We can thus make the following assumption on the latent samples  $\mathbf{h}_t$  without loss of generality.

**Assumption 2.** *The latent samples  $\mathbf{h}_t$  abide a standard normal distribution, i.e.,  $\mathbf{h}_t \sim \mathcal{N}(\cdot; 0, \mathbf{I})$ .*

**Remark 1.** *If the state transition matrix  $\mathbf{A}$  is given, and the process is stationary, i.e.,  $p(\mathbf{h}_t) = p(\mathbf{h}_{t+1})$  for all  $t$ , then Assumption 2 essentially identifies the process noise model. Namely, we have*

$$\begin{aligned} \mathbb{E}_{\mathbf{h}_{t+1}} [\mathbf{h}_{t+1}\mathbf{h}_{t+1}^\top] &= \mathbb{E}_{\mathbf{h}_t, \mathbf{v}_t} [(\mathbf{A}\mathbf{h}_t + \mathbf{v}_t)(\mathbf{A}\mathbf{h}_t + \mathbf{v}_t)^\top] \\ &= \mathbb{E}_{\mathbf{h}_t, \mathbf{v}_t} [\mathbf{A}\mathbf{h}_t\mathbf{h}_t^\top\mathbf{A}^\top + \mathbf{v}_t\mathbf{h}_t^\top\mathbf{A}^\top + \mathbf{A}\mathbf{h}_t\mathbf{v}_t^\top + \mathbf{v}_t\mathbf{v}_t^\top] \\ &= \mathbb{E}_{\mathbf{h}_t, \mathbf{v}_t} [\mathbf{A}\mathbf{h}_t\mathbf{h}_t^\top\mathbf{A}^\top] + \mathbb{E}_{\mathbf{v}_t} [\mathbf{v}_t\mathbf{v}_t^\top] \\ &= \mathbf{A}\mathbf{A}^\top + \mathbb{E}_{\mathbf{v}_t} [\mathbf{v}_t\mathbf{v}_t^\top]. \end{aligned} \quad (4)$$

In order to make sure that the latent states  $\mathbf{h}_t$  remain standard Gaussian in sequential synthesis scenarios, i.e.,  $\mathbb{E}_{\mathbf{h}_{t+1}} [\mathbf{h}_{t+1}\mathbf{h}_{t+1}^\top] = \mathbf{I}$ , we just need to ensure that the process noise is zero-mean and has the covariance matrix  $\mathbf{I} - \mathbf{A}\mathbf{A}^\top$ .  $\diamond$

Before we propose an algorithm to jointly learn  $\mathbf{C}$  and  $\mathbf{A}$  by means of deep generative architecture, we further characterize the problem in order to investigate its feasibility and motivate our approach. For a model Eq. (3) to substitute Eq. (2), it needs to transform the transition  $\varphi$  to a multiplication by a matrix and the perturbation  $\psi_{\mathbf{v}_t}$  to a superposition by zero-mean Gaussian noise. The latter is generally easy to achieve, if the perturbation is sufficiently small and  $\mathbf{C}$  is a diffeomorphism. In the remainder of the section, we thus focus on linearization of  $\varphi$ .

### 3.2 LOCAL LINEARIZATION

A common linearization method in control theory is to approximate the dynamic system equation by a first-order Taylor polynomial around an equilibrium point (Perko (2013)). Clearly, for this to be possible, such a point needs to exist. We thus propose the following assumption on the transition function  $\varphi$  of Eq. (2).

**Assumption 3.** The differentiable self map  $\varphi: \mathbb{M} \rightarrow \mathbb{M}$  has at least one fixed point  $\mathbf{y}^*$ .

However, such an approach might fail for real-world processes due to the curse of dimensionality. A remedy is to employ an representation function  $\Gamma$  that maps the system observations to a lower-dimensional latent space prior to performing the linearization. We call such a function a *local linearizer*.

**Definition 1.** Let  $\mathbb{M} \subset \mathbb{R}^d$  be an  $n$ -dimensional manifold,  $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}^d$  be differentiable on  $\varphi(\mathbb{M}) = \mathbb{M}$ , and  $\Gamma: \mathbb{R}^d \rightarrow \mathbb{R}^n$  be a diffeomorphism on  $\mathbb{M}$ . The map  $\Gamma$  is said to be a local linearizer at  $\mathbf{y}^* \in \mathbb{M}$  of  $\varphi$ , if there exists a matrix  $\Phi \in \mathbb{R}^{n \times n}$  such that the following equality holds true with  $\mathbf{y} \in \mathbb{R}^d$ .

$$\lim_{\|\mathbf{y} - \mathbf{y}^*\| \rightarrow 0} \frac{\|\varphi(\mathbf{y}) - \Gamma^{-1}(\Phi \Gamma(\mathbf{y}))\|}{\|\mathbf{y} - \mathbf{y}^*\|} = 0. \quad (5)$$

Generally speaking, local linearization is made possible by moving the fixed point to the origin of a new coordinate system. More precisely, the following proposition holds.

**Proposition 1.** Let  $\mathbb{M} \subset \mathbb{R}^d$  be an  $n$ -dimensional manifold. Furthermore, let  $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}^d$  be differentiable on  $\varphi(\mathbb{M}) = \mathbb{M}$ , and  $\Gamma: \mathbb{R}^d \rightarrow \mathbb{R}^n$  be a diffeomorphism on  $\mathbb{M}$ . If  $\mathbf{y}^* \in \mathbb{M}$  is a fixed point of  $\varphi$ , then the following map,

$$\begin{aligned} \Gamma' : \mathbb{R}^n &\rightarrow \mathbb{M} \\ \mathbf{y} &\mapsto \Gamma(\mathbf{y}) - \Gamma(\mathbf{y}^*), \end{aligned} \quad (6)$$

locally linearizes  $\varphi$ , if the rows and columns of the Jacobi matrix  $\mathbf{J}_\varphi$  of  $\varphi$  at  $\mathbf{y}^*$  lie in the row space of the Jacobi matrix  $\mathbf{J}_\Gamma$  of  $\Gamma$  at  $\mathbf{y}^*$ .

*Proof.* See Appendix A.1 □

The Jacobi matrix property is only needed for consistency with Definition 1, where the limit  $\mathbf{y} \rightarrow \mathbf{y}^*$  is approached from any possible direction in  $\mathbb{R}^d$ . It can be ignored, if we restrict the analysis entirely to the manifold  $\mathbb{M}$ . Proposition 1 states that if  $\varphi$  is a differentiable self map on  $\mathbb{M}$  with a fixed point  $\mathbf{y}^* \in \mathbb{M}$ , and a diffeomorphism from  $\mathbb{M}$  to  $\mathbb{R}^n$  exists, then there is a representation in which  $\varphi$  can be approximated by a linear function for points on  $\mathbb{M}$  that are not too far away from  $\mathbf{y}^*$ .

Note that even if the Jacobi matrix  $\mathbf{J}_\varphi$  can be sufficiently well estimated from the data, the linear approximation directly in the observation space  $\mathbb{R}^d$  will likely be worse than via reparametrization of a local linearizer  $\Gamma'$  that functions as a chart of  $\mathbb{M}$ . The reason is that predictions directly via  $\mathbf{J}_\varphi$  can not be assumed to remain on  $\mathbb{M}$ . Figuratively speaking, the aim of the local linearizer is to bend the space spanned by the rows and columns of  $\mathbf{J}_\varphi$  to match the shape of  $\mathbb{M}$ .

### 3.3 SEPARATE LEARNING

Local linearizability is a useful concept to analyze the feasibility of the problem at hand. However, it does not take into account that two local linearizers can have significantly different linearization properties on a global scale. Moreover, it does not provide any instructions on how to find an appropriate representation  $\Gamma$  and matrix  $\Phi$ . Traditionally, this task is approached by learning  $\Gamma$  and  $\Phi$  separately while considering the sampled observations of the process globally, as will be described in the following.

To characterize the problem of global linearization, we need to introduce a measure. It is sensible to consider the expected squared distance between the result of a transformation  $\varphi$  and its linear approximation. Since it is not possible to have an analytic expression for the distribution on the data manifold  $\mathbb{M}$ , the latent space is often considered as a heuristic. Because it is a common model assumption for deep generative models and in accordance with Assumption 2, we can assume standard Gaussian distribution for the latent space.

Let  $\Gamma: \mathbb{M} \rightarrow \mathbb{R}^n$  denote a data representation mapping and  $\varphi: \mathbb{M} \rightarrow \mathbb{M}$  the transition function to be linearized. Consider the following expression

$$q(\Gamma, \varphi, \Phi) = \frac{\mathbb{E}_{\mathbf{y}} [\|\Phi \Gamma(\mathbf{y}) - \Gamma(\varphi(\mathbf{y}))\|^2]}{2n}. \quad (7)$$

The denominator  $2n$  is a normalization factor. If  $\Gamma(\mathbf{y})$  has standard Gaussian distribution, then  $q(\Gamma, \varphi, \Phi) < 1$  denotes that the linear prediction with  $\Phi$  is smaller than the expected distance of two independently drawn samples. A common method of linearizing  $\varphi$ , is to first choose a representation  $\Gamma$  that is assumed to have good linearization properties, and then inferring  $\Phi$  by minimizing  $Eq. (7)$  (Doretto et al. (2003); Chan & Vasconcelos (2007); Han et al. (2017)). In that case, the solution is given by

$$\hat{\Phi} = \mathbb{E}_{\mathbf{y}}[\Gamma(\varphi(\mathbf{y}))\Gamma(\mathbf{y})^\top], \quad (8)$$

if  $\Gamma(\mathbf{y})$  has standard Gaussian distribution.

However, the drawback of this approach is the difficulty to find an appropriate model for  $\Gamma$  that can be assumed to linearize the transformation  $\varphi$ . Moreover, separate learning implicitly assumes that a small linearization error in the embedded space will carry over to a small error in the observation space, but such an assumption is hard to justify given the high dimension of the problem. We therefore propose to approach the problem by learning the representation and the linear transition jointly, by approximating the data distribution directly via deep generative models.

## 4 JOINT LEARNING VIA DEEP GENERATIVE MODELS

### 4.1 MARKOV ASSUMPTIONS

We now turn to the problem of finding a model Eq. (3) that approximates Eq. (2). In the following, we will provide the preliminaries to lay out a training procedure to infer  $C$  and  $\mathbf{A}$  from observed sequences of visual processes. A sequence of length  $N$  of a visual process  $\mathcal{Y}$  is viewed as a realization of the random variable  $\mathbf{y}^N = [\mathbf{y}_1, \dots, \mathbf{y}_N]$ . The according sequence in the latent space is a realization of the random variable  $\mathbf{h}^N = [\mathbf{h}_1, \dots, \mathbf{h}_N]$ . Consider the random variable

$$\tilde{\mathbf{y}}^N = [\tilde{\mathbf{y}}_1, \dots, \tilde{\mathbf{y}}_N] := [C(\mathbf{h}_1), \dots, C(\mathbf{h}_N)]. \quad (9)$$

In order to model  $\mathcal{Y}$  by Eq. (3), we need to make sure that the probability distributions of  $\mathbf{y}^N$  and  $\tilde{\mathbf{y}}^N$  coincide for any  $N \in \mathbb{N}$ . By taking into account Assumption 1, this is equivalent to demanding that the joint probability of two succeeding frames coincide.

The joint probability distribution of  $\mathbf{h}_t$  and  $\mathbf{h}_{t+1}$  is zero-mean Gaussian and, more specifically,

$$\begin{bmatrix} \mathbf{h}_t \\ \mathbf{h}_{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \end{bmatrix}; 0, \begin{bmatrix} \mathbf{I} & \mathbf{A}^\top \\ \mathbf{A} & \mathbf{I} \end{bmatrix}\right), \quad (10)$$

holds due to Assumption 2. To summarize, we are looking for the function  $C$  and a matrix  $\mathbf{A}$ , such that the random variable

$$\tilde{\mathbf{y}}^2 = [C(\mathbf{h}_1) \quad C(\mathbf{h}_2)], \quad (11)$$

has the same probability distribution as  $\mathbf{y}^2$ .

### 4.2 DEEP GENERATIVE MODELS

Estimation of probability distributions for high-dimensional data is still an ongoing research topic in deep learning. The problem is typically framed as an approximation of a function  $f_\theta$ , parameterized by a vector  $\theta$  in a finite-dimensional euclidean space. The purpose of  $f_\theta$  is to map realizations of low-dimensional, standard Gaussian noise to samples that abide the data distribution. Typically, the function  $f_\theta$  is realized by a neural network with trainable weights  $\theta$ .

Most of the widely employed deep generative models, including the GAN, and the VAE, adopt this approach. What these models differ in, is merely the way the network parameters  $\theta$  are trained. For this work, we employ the Wasserstein GAN and the VAE to evaluate the proposed method. Due to space constraints, we will not review these techniques here and refer the reader to Arjovsky et al. (2017) for the Wasserstein GAN and to Doersch (2016) for the VAE. For the following sections, it is sufficient to assume that these techniques are capable of approximating a function  $f_\theta$  that maps from a low-dimensional standard Gaussian to a high-dimensional data distribution.

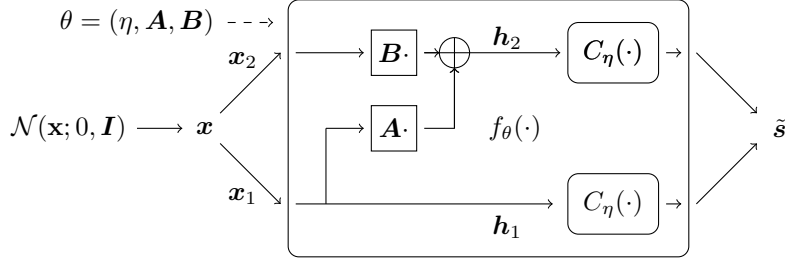


Figure 1: Generative network with dynamic layer. The dashed arrow indicates the trainable parameters

#### 4.3 THE DYNAMIC LAYER

We now focus on how to learn a matrix  $\mathbf{A}$  and a function  $C$  from a finite sequence  $\mathbf{y}^N \in \mathbb{R}^{d \times N}$  of a visual process such that it can be described by Eq. (3). To this end, recall that, as discussed in Section 4.1, we only need to consider the joint probability of two succeeding observed frames. The first step is thus to generate a training set of frame pairs from  $\mathbf{y}^N$  as  $\{\mathbf{s}_1, \dots, \mathbf{s}_{N-1}\}$ , where

$$\mathbf{s}_t = [\mathbf{y}_t^\top \quad \mathbf{y}_{t+1}^\top]^\top \quad (12)$$

denotes a vectorized pair of frames. We treat the samples as realizations of the random variable  $\mathbf{s}$  and are looking for an architecture to learn the function  $C$  and a matrix  $\mathbf{A}$  such that the probability distribution of the random variable

$$[C(\mathbf{h}_1)^\top \quad C(\mathbf{h}_2)^\top]^\top, \quad (13)$$

where  $\mathbf{h}_1, \mathbf{h}_2$  have the joint distribution described by Eq. (10), coincides with the probability distribution of  $\mathbf{s}$ , i.e.

$$p(\mathbf{s}) = p\left([C(\mathbf{h}_1)^\top \quad C(\mathbf{h}_2)^\top]^\top\right). \quad (14)$$

It turns out that we can use a deep generative architecture to accomplish this task. Remember, that a function  $f_\theta$  learned by a deep generative model is capable of transforming samples  $\mathbf{x}$  of standard Gaussian noise to samples from a high dimensional data distribution. Let  $\theta = (\mathbf{A}, \mathbf{B}, \eta)$  contain the matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$  such that the constraint

$$\mathbf{A}\mathbf{A}^\top + \mathbf{B}\mathbf{B}^\top = \mathbf{I}_n \quad (15)$$

is fulfilled, and the parameter  $\eta$  that defines the model for  $C_\eta = C$ . Consider the following definition for  $f_\theta$ .

$$f_\theta : \mathbb{R}^{2n} \rightarrow \mathbb{R}^{2d}, \quad \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \mapsto \begin{bmatrix} C_\eta(\mathbf{x}_1) \\ C_\eta(\mathbf{A}\mathbf{x}_1 + \mathbf{B}\mathbf{x}_2) \end{bmatrix}. \quad (16)$$

Assume that we can train  $f_\theta$ , as defined in Eq. (16), to map  $\mathbf{x} \sim \mathcal{N}(\mathbf{x}; 0, \mathbf{I}_{2n})$  to  $\tilde{\mathbf{s}} = f_\theta(\mathbf{x})$  with the probability distribution

$$\tilde{\mathbf{s}} \sim p(\mathbf{s}). \quad (17)$$

Then,  $\mathbf{A}$  and  $C = C_\eta$  indeed fulfill the condition in Eq. (14) for  $\mathbf{h}_1, \mathbf{h}_2$  with joint probability distribution Eq. (10). In order to implement  $f_\theta$  by a neural network, we propose the architecture depicted in Fig. 1. The first layer is linear and will be referred to as the *dynamic layer* in the following. The dynamic layer implements the multiplication with the matrix

$$\mathbf{F} = \begin{bmatrix} \mathbf{I}_n & \mathbf{B} \\ \mathbf{A} & \mathbf{B} \end{bmatrix}. \quad (18)$$

The output of the dynamic layer is split into an upper half  $\mathbf{h}_1$  and a lower half  $\mathbf{h}_2$  and both halves are fed to the subnetwork that implements the observer function  $C_\eta$ . The weights of the dynamic layer contain the matrices  $\mathbf{A}, \mathbf{B}$ . Thus, they can be trained along with  $\eta$ , by back-propagation of the stochastic gradient. However, we need to make sure that the stationarity constraint in Eq. (15) is not violated. This can be done by adding a regularizer to the loss function  $L$  of the Deep Generative Network. We thus substitute  $L(\theta)$  by

$$\tilde{L}(\theta) = L(\theta) + \lambda \|\mathbf{A}\mathbf{A}^\top + \mathbf{B}\mathbf{B}^\top - \mathbf{I}_n\|_F^2, \quad (19)$$

where  $\lambda > 0$  is a regularizer constant. Here,  $L$  refers to the variational bound for the VAE or to the Generator loss for the (Wasserstein) GAN.

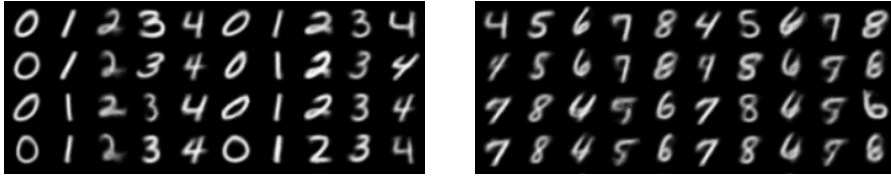


Figure 2: Synthesis result for sequences of MNIST digits with VAE

## 5 EXPERIMENTS

### 5.1 OVERVIEW

Due to its impressive results on natural images, we employ the generator of the Deep Convolutional GAN (DCGAN, Radford et al. (2015)) as the foundation for implementing the observer function  $C_\eta$ . We use it in combination with an affine layer at the input. This layer serves three purposes. Firstly, it accounts for changes of bases in the latent space in order to conform with Assumption 2. Secondly, it includes a bias so that the fixed point of the transition function can be transformed to the origin according to Proposition 1. Finally, it further reduces the latent dimension to make the search for the transition matrix  $A$  feasible.

We train  $C_\eta$  and  $A$  simultaneously by means of the dynamic layer introduced in Section 4.3. To this end, we integrate the construction in Fig. 1 as the decoder of a VAE (without batch normalization) and the generator of a Wasserstein GAN (with batch normalization), respectively. As the critic (discriminator) network of the Wasserstein GAN, we use the original discriminator of the DCGAN but double the number of output channels in order to match the dimensions of the training frame pairs. As the encoder of the VAE, we use the discriminator of the DCGAN, but adapt the number of output channels to be  $2n$ , where  $n$  is the latent dimension of the model. The latent dimension is set to  $n = 10$  for all experiments.

Synthesis is performed by sampling from the MAR model described by

$$\mathbf{h}_{t+1} = A\mathbf{h}_t + B\mathbf{v}_t, \quad \mathbf{v}_t \sim \text{i.i.d. } \mathcal{N}(\mathbf{v}_t, 0, \mathbf{I}_n), \quad (20)$$

and mapping the latent states to the observation space by means of  $C_\eta$ . The initial latent state  $\mathbf{h}_0$  is sampled from Gaussian white noise for the Wasserstein GAN experiments. For the VAE experiments, it is estimated by applying the encoder to an observation frame pair from the training sequence. PyTorch code for reproducing the experiments will be made available online upon publication, along with the entire set of results.

### 5.2 VARIATIONAL AUTOENCODER

An isotropic standard Gaussian model  $\mathcal{N}(\mathbf{s}; f_\theta(\mathbf{x}), \sigma^2)$  is assumed for the conditional data distribution. The Adam optimizer with a step size of  $2.5 * 10^{-4}$  is used to train the architecture. The regularizer constant is set to  $\lambda = 100$ . We carry out two series of experiments.

In the first series of experiments, we use the *MNIST* dataset to generate repeating sequences of hand written digits. We choose  $\sigma = 4$  and use sequences of length 10000 for training the network for 25 epochs. Most of the number sequences we test can be well reproduced. Fig. 2 depicts the synthesis results for the sequences 0123401234... and 4567845678.... Occasional jumps occur, due to the stochastic nature of the model. An observation we make is that higher values of  $\sigma$  improve the probability of synthesizing the correct sequence, but decrease the variability of digit shapes.

To investigate the linearization of geometrical transformations, we employ Category 0 of the *Small NORB* dataset that contains pictures of miniature animals under 6 different lighting conditions, 9 elevational and 18 azimuthal poses. We train our architecture for 100 epochs with  $\sigma = \sqrt{5}$  to synthesize sequences depicting azimuthal rotations of  $20^\circ$  per frame. Generally, the synthesis provides good results. For each one of the five animals, the synthesis yields clearly recognizable rotation movements. However, the algorithm tends to confuse pairs of opposite poses, as can be observed in Fig. 3a, where the poses are flipped after the fourth frame in each row. With regard to Section 3.2,

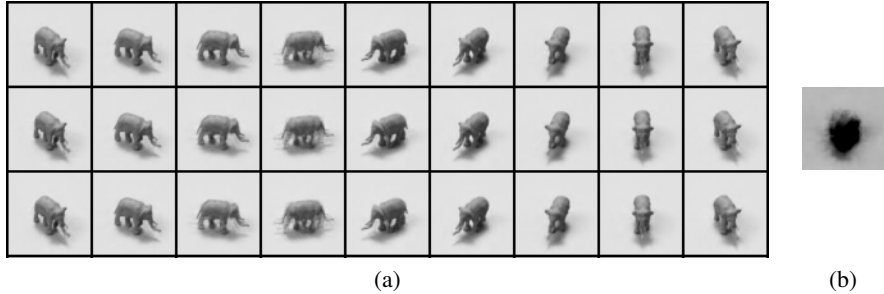


Figure 3: Synthesis result for NORB sequences (a) and learned fixed point (b)

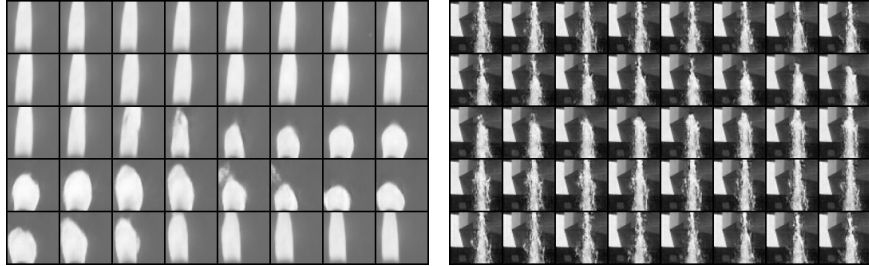


Figure 4: Synthesis result for sequences of the UCLA-50 dataset with Wasserstein GAN

two explanations for this can be identified. On the one hand, it is possible that the assumptions of Proposition 1 are not fulfilled, and no diffeomorphism, i.e. a bijective mapping exists from the image manifold to the embedded space. In that case, the best we can do is finding a non-injective local diffeomorphism that assigns more than one point on the image manifold to a point in the embedded space. On the other hand, it is thinkable, that an actual diffeomorphism exists, but that the euclidean distance implicitly minimized by the VAE is a bad estimation of the geodesic distance on the manifold. Fig. 3b depicts the learned fixed point of the transition function. It can be thought of as a rotationally invariant structure under limited exposure of light from above.

### 5.3 WASSERSTEIN GENERATIVE ADVERSARIAL NETWORK

We employ the Wasserstein GAN for dynamic texture synthesis experiments and set the regularizer to  $\lambda = 1$ . We firstly test our method on the cropped *UCLA-50* dataset of grayscale DTs. The dataset contains 50 classes of 4 sequences of length 75 each. We use one class at a time for the experiments. The architecture is trained for 2000 epochs via RMSPROP with step size  $2.5 * 10^{-4}$ . Fig. 4 depicts the results for the classes `candle` and `fountain-c-far`.

Finally, we evaluate the method on three RGB sequences. We run the optimization via RMSPROP with step size  $5 * 10^{-5}$  for 900 epochs. Fig. 5 depicts the result for the `firepot`, `springwater` and `waterfountain` sequences taken from Xie et al. (2017). Despite the low complexity of our model, we believe that the difference in quality compared to the state of the art results reported in Xie et al. (2017) is negligible.

## 6 CONCLUSION

This work presents an approach to learn embedded MAR models from image sequences. We motivate the feasibility of this approach by introducing the concept of local linearizability and propose a joint learning procedure that employs deep generative models in combination with an additional linear component, the dynamic layer. We report first positive results on low-resolution visual processes, where a first-order Markov property can be assumed, and hope to shed some light on the nature of linearization. A possible future research direction is improving the theoretical understanding of linearizing representations and their applicability outside of stationary visual processes.





Figure 5: Synthesis result for RGB sequences. Frames are rescaled to match the original aspect ratio

## REFERENCES

- Bijan Afsari and René Vidal. The alignment distance on spaces of linear dynamical systems. In *Decision and Control (CDC), 2013 IEEE 52nd Annual Conference on*, pp. 1162–1167. IEEE, 2013.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017.
- Rakesh Chalasani and Jose C Principe. Deep predictive coding networks. *arXiv preprint arXiv:1301.3541*, 2013.
- Antoni B Chan and Nuno Vasconcelos. Classifying video with kernel dynamic textures. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1–6. IEEE, 2007.
- Taco Cohen and Max Welling. Group equivariant convolutional networks. In *International Conference on Machine Learning*, pp. 2990–2999, 2016.
- Taco S Cohen and Max Welling. Transformation properties of learned visual representations. *arXiv preprint arXiv:1412.7659*, 2014.
- CR De Souza, A Gaidon, Y Cabon, and AM Lopez Pena. Procedural generation of videos to train deep action recognition networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- Carl Doersch. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*, 2016.
- Gianfranco Doretto, Alessandro Chiuso, Ying Nian Wu, and Stefano Soatto. Dynamic textures. *International Journal of Computer Vision*, 51(2):91–109, 2003.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pp. 2672–2680, 2014.
- Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT press Cambridge, 2016.
- Ross Goroshin, Michael F Mathieu, and Yann LeCun. Learning to linearize under uncertainty. In *Advances in Neural Information Processing Systems*, pp. 1234–1242, 2015.
- Tian Han, Yang Lu, Song-Chun Zhu, and Ying Nian Wu. Alternating back-propagation for generator network. In *AAAI*, volume 3, pp. 13, 2017.
- Geoffrey E Hinton, Alex Krizhevsky, and Sida D Wang. Transforming auto-encoders. In *International Conference on Artificial Neural Networks*, pp. 44–51. Springer, 2011.
- Matthew Johnson, David K Duvenaud, Alex Wiltschko, Ryan P Adams, and Sandeep R Datta. Composing graphical models with neural networks for structured representations and fast inference. In *Advances in neural information processing systems*, pp. 2946–2954, 2016.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Rahul G Krishnan, Uri Shalit, and David Sontag. Deep kalman filters. *arXiv preprint arXiv:1511.05121*, 2015.
- Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *Advances in Neural Information Processing Systems*, pp. 700–708, 2017a.
- Wenqian Liu, Abhishek Sharma, Octavia Camps, and Mario Sznaier. Dyan: A dynamical atoms-based network for video prediction. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 170–185, 2018.
- Ziwei Liu, Raymond Yeh, Xiaoou Tang, Yiming Liu, and Aseem Agarwala. Video frame synthesis using deep voxel flow. In *International Conference on Computer Vision (ICCV)*, volume 2, 2017b.

- Michael Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond mean square error. *arXiv preprint arXiv:1511.05440*, 2015.
- Roland Memisevic. Learning to relate images. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1829–1846, 2013.
- Lawrence Perko. *Differential equations and dynamical systems*, volume 7. Springer Science & Business Media, 2013.
- Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. 2015.
- Nitish Srivastava, Elman Mansimov, and Ruslan Salakhudinov. Unsupervised learning of video representations using lstms. In *International conference on machine learning*, pp. 843–852, 2015.
- Carl Vondrick and Antonio Torralba. Generating the future with adversarial transformers. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. In *Advances In Neural Information Processing Systems*, pp. 613–621, 2016.
- Manuel Watter, Jost Springenberg, Joschka Boedecker, and Martin Riedmiller. Embed to control: A locally linear latent dynamics model for control from raw images. In *Advances in neural information processing systems*, pp. 2746–2754, 2015.
- Xian Wei, Yuanxiang Li, Hao Shen, Fang Chen, Martin Kleinsteuber, and Zhongfeng Wang. Dynamical textures modeling via joint video dictionary learning. *IEEE Transactions on Image Processing*, 26(6):2929–2943, 2017.
- Jianwen Xie, Song-Chun Zhu, and Ying Nian Wu. Synthesizing dynamic patterns by spatial-temporal generative convnet. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7093–7101, 2017.
- Tianfan Xue, Jiajun Wu, Katherine Bouman, and Bill Freeman. Visual dynamics: Probabilistic future frame synthesis via cross convolutional networks. In *Advances in Neural Information Processing Systems*, pp. 91–99, 2016.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. *arXiv preprint*, 2017.

## A APPENDIX

### A.1 PROOF OF PROPOSITION 1

Let  $\Phi \in \mathbb{R}^{n \times n}$  be a matrix. Since  $\Gamma'$  is a diffeomorphism, we define  $\phi = \Gamma'^{-1} \circ \Phi \circ \Gamma'$ . We need to show that

$$\lim_{\|\mathbf{y} - \mathbf{y}^*\| \rightarrow 0} \frac{\|\varphi(\mathbf{y}) - \phi(\mathbf{y})\|}{\|\mathbf{y} - \mathbf{y}^*\|} = 0 \quad (21)$$

holds. Let us denote the Jacobian of  $\phi$  at  $\mathbf{y}^*$  by  $\mathbf{J}_\phi$ . Because  $\Gamma'$  maps  $\mathbf{y}^*$  to the origin, we can reformulate the requirement as

$$\begin{aligned} & \lim_{\|\mathbf{y} - \mathbf{y}^*\| \rightarrow 0} \frac{\|\varphi(\mathbf{y}) - \mathbf{y}^* - \phi(\mathbf{y}) + \mathbf{y}^*\|}{\|\mathbf{y} - \mathbf{y}^*\|} \\ &= \lim_{\|\mathbf{y} - \mathbf{y}^*\| \rightarrow 0} \frac{\|\varphi(\mathbf{y}) - \varphi(\mathbf{y}^*) - \phi(\mathbf{y}) + \phi(\mathbf{y}^*)\|}{\|\mathbf{y} - \mathbf{y}^*\|} \\ &= \lim_{\|\mathbf{y} - \mathbf{y}^*\| \rightarrow 0} \frac{\|\varphi(\mathbf{y}) - \varphi(\mathbf{y}^*) - \mathbf{J}_\phi(\mathbf{y} - \mathbf{y}^*)\|}{\|\mathbf{y} - \mathbf{y}^*\|} = 0. \end{aligned} \quad (22)$$

This requirement is fulfilled if the Jacobi matrices of  $\varphi$  and  $\phi$  coincide, i.e.

$$\mathbf{J}_\phi = \mathbf{J}_\varphi. \quad (23)$$

The Jacobian of  $\phi$  at  $\mathbf{y}^*$  is given, according to the chain rule, by

$$\mathbf{J}_\phi = \mathbf{J}_{\Gamma^{-1}} \mathbf{\Phi} \mathbf{J}_\Gamma, \quad (24)$$

where  $\mathbf{J}_{\Gamma^{-1}} \in \mathbb{R}^{d \times n}$  is the Jacobian matrix of  $\Gamma^{-1}$  at  $\Gamma(\mathbf{y}^*)$ . A matrix  $\mathbf{\Phi} \in \mathbb{R}^{n \times n}$  can be always found, such that Eq. (23) is fulfilled, if the columns and rows of  $\mathbf{J}_\phi$  lie in the column space of  $\mathbf{J}_{\Gamma^{-1}}$  and the row space of  $\mathbf{J}_\Gamma$ , respectively. The column space of  $\mathbf{J}_{\Gamma^{-1}}$  coincides with the row space of  $\mathbf{J}_\Gamma$ , due to the identity  $\mathbf{J}_\Gamma \mathbf{J}_{\Gamma^{-1}} = \mathbf{I}_n$ . The statement of the proposition follows.

## A.2 GENERALIZATION TO HIGHER-ORDER MARKOV PROCESSES

The proposed method is designed to linearize first-order Markov processes only. This is in line with our empirical evaluation, in which the method was not capable of synthesizing sequences with persisting motions of non-constant velocities, e.g. the videos of running animals in Xie et al. (2017). However, the presented approach can be theoretically generalized to Markov processes of order  $m \geq 1$ . We present a procedure that could be employed for this purpose. Further investigation is necessary to evaluate the feasibility and practical applicability of such an approach.

The general procedure is as follows.

1. Set up an appropriate block matrix model for the joint probability of  $m + 1$  succeeding observations.
2. Build a dynamic layer that performs a multiplication with a lower-triangular  $(m + 1) \times (m + 1)$  block matrix  $\mathbf{F}$ .
3. Introduce regularizers to preserve the block Toeplitz structure of the covariance matrix,
4. Compute the MAR parameters from the resulting covariance matrix.

We illustrate the procedure for the case  $m = 2$ .

1. Due to the stationarity assumption, the covariance matrix for three succeeding observations  $\mathbf{h}_t, \mathbf{h}_{t+1}, \mathbf{h}_{t+2}$  has the form

$$\mathbf{\Sigma}_2 = \begin{bmatrix} \mathbf{I}_n & \text{Cov}(\mathbf{h}_t, \mathbf{h}_{t+1}) & \text{Cov}(\mathbf{h}_t, \mathbf{h}_{t+2}) \\ \text{Cov}(\mathbf{h}_t, \mathbf{h}_{t+1})^\top & \mathbf{I}_n & \text{Cov}(\mathbf{h}_{t+1}, \mathbf{h}_{t+2}) \\ \text{Cov}(\mathbf{h}_t, \mathbf{h}_{t+2})^\top & \text{Cov}(\mathbf{h}_{t+1}, \mathbf{h}_{t+2})^\top & \mathbf{I}_n \end{bmatrix}. \quad (25)$$

2. The matrix describing the dynamic layer has the form

$$\mathbf{F} = \begin{bmatrix} \mathbf{I}_n & & \\ \mathbf{F}_1 & \mathbf{F}_2 & \\ \mathbf{F}_3 & \mathbf{F}_4 & \mathbf{F}_5 \end{bmatrix}. \quad (26)$$

The outputs of the dynamic layer will thus have the distribution

$$p\left(\begin{bmatrix} \mathbf{h}_1^\top & \mathbf{h}_2^\top & \mathbf{h}_3^\top \end{bmatrix}^\top\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{h}_1^\top & \mathbf{h}_2^\top & \mathbf{h}_3^\top \end{bmatrix}^\top; 0, \mathbf{F} \mathbf{F}^\top\right). \quad (27)$$

We thus need to achieve

$$\mathbf{\Sigma}_2 \approx \mathbf{F} \mathbf{F}^\top = \begin{bmatrix} \mathbf{I}_n & \mathbf{F}_1^\top & \mathbf{F}_3^\top \\ \mathbf{F}_1 & \mathbf{F}_1 \mathbf{F}_1^\top + \mathbf{F}_2 \mathbf{F}_2^\top & \mathbf{F}_1 \mathbf{F}_3^\top + \mathbf{F}_2 \mathbf{F}_4^\top \\ \mathbf{F}_3 & \mathbf{F}_3 \mathbf{F}_1^\top + \mathbf{F}_4 \mathbf{F}_2^\top & \mathbf{F}_3 \mathbf{F}_3^\top + \mathbf{F}_4 \mathbf{F}_4^\top + \mathbf{F}_5 \mathbf{F}_5^\top \end{bmatrix}. \quad (28)$$

3. This can be ensured by using the regularizer

$$R(\mathbf{F}) = \lambda_1 \|\mathbf{F}_1 \mathbf{F}_1^\top + \mathbf{F}_2 \mathbf{F}_2^\top - \mathbf{I}_n\|_F^2 + \lambda_2 \|\mathbf{F}_3 \mathbf{F}_3^\top + \mathbf{F}_4 \mathbf{F}_4^\top + \mathbf{F}_5 \mathbf{F}_5^\top - \mathbf{I}_n\|_F^2 + \lambda_3 \|\mathbf{F}_3 \mathbf{F}_1^\top + \mathbf{F}_4 \mathbf{F}_2^\top - \mathbf{F}_1 \mathbf{F}_3^\top\|_F^2, \quad (29)$$

with  $\lambda_1, \lambda_2, \lambda_3 > 0$ .

4. A second-order MAR model has the form

$$\mathbf{h}_{t+2} = \mathbf{A}_0 \mathbf{h}_t + \mathbf{A}_1 \mathbf{h}_{t+1} + \mathbf{B} \mathbf{v}_t, \quad \mathbf{v}_t \sim \mathcal{N}(\mathbf{v}_t, 0, \mathbf{I}_n). \quad (30)$$

By our assumptions on the process, this yields the system of equations,

$$\begin{aligned}
\text{Cov}(\mathbf{h}_t, \mathbf{h}_t) &= \mathbf{A}_0 \mathbf{A}_0^\top + \mathbf{A}_1 \mathbf{A}_1^\top + \mathbf{A}_0 \text{Cov}(\mathbf{h}_t, \mathbf{h}_{t+1}) \mathbf{A}_1^\top \\
&\quad + \mathbf{A}_1 \text{Cov}(\mathbf{h}_t, \mathbf{h}_{t+1})^\top \mathbf{A}_0^\top + \mathbf{B} \mathbf{B}^\top \\
\text{Cov}(\mathbf{h}_t, \mathbf{h}_{t+1}) &= \text{Cov}(\mathbf{h}_t, \mathbf{h}_{t+1})^\top \mathbf{A}_0^\top + \mathbf{A}_1^\top, \\
\text{Cov}(\mathbf{h}_t, \mathbf{h}_{t+2}) &= \mathbf{A}_0^\top + \text{Cov}(\mathbf{h}_t, \mathbf{h}_{t+1}) \mathbf{A}_1^\top.
\end{aligned} \tag{31}$$

Thus, we can infer  $\mathbf{A}_0$ ,  $\mathbf{A}_1$  and  $\mathbf{B}$  via solving

$$\begin{aligned}
\mathbf{F}_1^\top &= \mathbf{F}_1 \mathbf{A}_0^\top + \mathbf{A}_1^\top, \\
\mathbf{F}_3 &= \mathbf{A}_0 + \mathbf{A}_1 \mathbf{F}_1, \\
\mathbf{I}_n &= \mathbf{A}_0 \mathbf{A}_0^\top + \mathbf{A}_1 \mathbf{A}_1^\top + \mathbf{A}_0 \mathbf{F}_1^\top \mathbf{A}_1^\top + \mathbf{A}_1 \mathbf{F}_1 \mathbf{A}_0^\top + \mathbf{B} \mathbf{B}^\top.
\end{aligned} \tag{32}$$