A Practical Approach to Causal Inference over Time

Anonymous submission

Abstract

In this paper, we focus on estimating the causal effect of an intervention over time on a dynamical system. To that end, we formally define causal interventions and their effects over time on discrete-time stochastic processes (DSPs). Then, we show under which conditions the equilibrium states of a DSP, both before and after a causal intervention, can be captured by a structural causal model (SCM). With such an equivalence at hand, we provide an explicit mapping from vector autoregressive models (VARs), broadly applied in econometrics, to linear, but potentially cyclic and/or affected by unmeasured confounders, SCMs. The resulting causal VAR framework allows us to perform causal inference over time from observational time series data. Our experiments on synthetic and real-world datasets show that the proposed framework achieves strong performance in terms of observational forecasting while enabling accurate estimation of the causal effect of interventions on dynamical systems. We demonstrate, through a case study, the potential practical questions that can be addressed using the proposed causal VAR framework.

1 Introduction

Dynamical systems often exhibit complex behaviors that unfold over time, leading to delayed responses and feedback loops. Importantly, understanding the causal effect of interventions within such systems is crucial across disciplines such as climate (Runge et al. 2019) and social sciences (Wunsch et al. 2022), where different time scales play a central role. For instance, monetary policy adjustments may have immediate effects on consumer spending, but their impact on inflation, employment, and economic growth only becomes evident in the medium-/long-term. Similarly, the consequences of human actions on climate change may take decades to manifest, with the risk of endorsing public policies that underestimate their relevance. To address these issues, it is essential to estimate the causal effect of interventions, or generally, to perform *causal inference over time*.

From the perspective of causality, Structural Causal Models (SCMs) provide a formal framework to perform causal inference from cross-sectional data. However, adapting existing methods to capture temporal dynamics remains a challenge (Bongers et al. 2021). Alternatively, temporal models, such as autoregressive models, offer practical methods for time-series analysis and forecasting (Lütkepohl 2005), but their formalization of causal effects is limited. First, they model interventions as *shocks* applied at a specific point in time, with effects that fade away after a certain period (Hyvärinen et al. 2010; Moneta et al. 2011). Second, they rely on Granger causality (Granger 1969) which is concerned with how well one variable can predict another rather than identifying causal relationships between them.

Our work combines the strengths of both frameworks, i.e., SCMs and autoregressive models, to enable robust reasoning about the causal effect of interventions on dynamical systems over time. To that end, we first introduce a formal definition of causal interventions on discrete-time stochastic processes (DSPs), proposing two alternatives, additive and forcing interventions. Second, we establish conditions under which the equilibrium state of a DSP can be represented by an SCM. Third, we develop a framework that maps VARs to linear SCMs, handling potentially cyclic structures and unmeasured confounders. Finally, our practical framework for causal inference over time from observational time-series data is empirically validated on synthetic and real-world datasets.

Related work The works most closely related to ours are these from Mooij, Janzing, and Schölkopf (2013) and Bongers, Blom, and Mooij (2018), as they theoretically connect dynamical systems to the causal semantics of SCMs via the equilibration of deterministic and random differential equations, and thus are capable of modeling cyclic causal mechanisms (Bongers et al. 2021). Our approach differs from this line of work in two key aspects: i) we focus on discrete-time dynamical systems parameterized using stochastic equations which, as stated by Bongers, Blom, and Mooij (2018), become particularly challenging for continuous-time processes; and ii) our mapping from autoregressive DSPs to SCMs provides not only a theoretical but also, to the best of our knowledge, the first data-driven framework for performing causal inference over time in dynamical systems.

2 Preliminaries and background

2.1 Structural Causal Models

Equation A SCM $\mathcal{M} = (\mathbf{F}, \mathbf{E})$ determines how a set of d endogenous (observed) random variables $\mathbf{X} := \{X^{(1)}, \ldots, X^{(d)}\}$ are obtained from a set of exogenous variables $\mathbf{E} := \{E_1, \ldots, E_d\}$, with prior distribution

 $p(\mathbf{E})$, via a set of structural equations $\mathbf{F} := \{X^{(i)} := f_i(\mathbf{PA}^{(i)}, \mathbf{E}^{(i)})\}_{i=1}^d$. Each f_i computes $X^{(i)}$ from its causal parents¹ $\mathbf{PA}^{(i)} \subseteq \mathbf{X}$ and a set $\mathbf{E}^{(i)} \subseteq \mathbf{E}$. We refer to \mathbf{X} as a solution of \mathcal{M} . We assume $\mathbf{PA}^{(i)}$ to be minimal, i.e., it only contains variables $X^{(j)}$ such that $\partial_{X^{(j)}} f_i \neq 0$. This formulation extends the definition in (Pearl 2009) to include cycles as in (Bongers, Blom, and Mooij 2018).

Graph A SCM \mathcal{M} induces a directed graph $\mathcal{G}_{\mathcal{M}} = (\mathcal{V}, \mathcal{E})$ that describes the functional dependencies in $\mathbf{F}: \mathcal{V}$ is the set of nodes for which V_i represents $X^{(i)}$ and \mathcal{E} is the set of the edges $(V_i, V_j) \in \mathcal{E} \iff X^{(i)} \in \mathbf{PA}^{(j)}$.

Intervention Besides describing the observational distribution $p(\mathbf{X})$, SCMs allow answering *interventional queries* about the effect of external manipulations, and enable *counterfactual queries* assessing what would have happened to a particular observation if one observed variable $X^{(i)}$ had taken a different value. An intervention \mathcal{I} on a SCM \mathcal{M} yields a new SCM $\mathcal{M}^{\mathcal{I}}$ for which one or more mechanisms $f_i(\mathbf{PA}^{(i)}, \mathbf{E}^{(i)})$ change to $\tilde{f}_i(\mathbf{PA}^{(i)}, \mathbf{\tilde{E}}^{(i)})$, where $\mathbf{P}\mathbf{\tilde{A}}^{(i)} \subseteq \mathbf{PA}^{(i)}$ and $\mathbf{\tilde{E}}^{(i)} \subseteq \mathbf{E}^{(i)}$. We refer to a *hard intervention* when f_i is replaced by a constant value $\alpha^{(i)}$, and $\mathbf{P}\mathbf{\tilde{A}}^{(i)} = \mathbf{\tilde{E}}^{(i)} = \emptyset$. This type of intervention is denoted by the do-operator $do(X^{(i)} = \alpha^{(i)})$. On the other hand, we refer to a *soft intervention* when at least one argument of f_i is retained. The causal effect CE of an intervention is evaluated in terms of differences between the values of the observable variables before and after the intervention \mathcal{I} , i.e.,

$$CE^{\mathcal{I}} = \mathbb{E}[X^{\mathcal{I}} - X].$$
(1)

2.2 Discrete-time Stochastic Processes

A discrete-time (vector) stochastic process (DSP) is a function $X : T \times \Omega \to \mathbb{R}^d$ where $t \in T$ is a time index in \mathbb{Z} , such that X_t (which denotes $X(t, \cdot)$) is a random variable on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$. We refer to $X(\omega)$ (which denotes $X(\cdot, \omega)$) as a *realization* or *trajectory* of X and denote the i-th component of X with $X^{(i)}$. Every DSP can be described through a difference equation (DE), i.e., a recurrence relation that allows computing X_t based on its past values. DEs can be categorized into three types (Bongers, Blom, and Mooij 2018): ordinary difference equations (ODE) describing deterministic processes; random difference equations (RDE), which involve randomness in the initial state X_0 and in the evolution parameters (see App. A); and stochastic difference equations, which describe inherently stochastic trajectories.

Equation A stochastic difference equation (SDE) describes a DSP via a functional relationship of the form

$$\boldsymbol{X}_{t} = \boldsymbol{f}(\boldsymbol{X}_{< t}) + \boldsymbol{g}(\boldsymbol{X}_{< t}) \odot \boldsymbol{\varepsilon}_{t}, \qquad (2)$$

where $X_{<t} := \{X_{t-1}, X_{t-2}, ...\}$ represents the trajectory up to time t, f represents the system's deterministic

mechanism, and the Hadamard product $g \odot \varepsilon_t$ is the inhomogeneous stochastic part, where ε_t denotes white noise, i.e., $\forall t, t' \in T$ $\mathbb{E}[\varepsilon_t] = \mathbf{0}, \mathbb{E}[\varepsilon_t^2] = \Sigma_{\varepsilon}^2, \mathbb{E}[\varepsilon_t \varepsilon_{t'}] = \mathbf{0}$. While for ODEs and RDEs trajectories $X(\omega)$ may asymptotically converge to an equilibrium, SDEs cannot exhibit such convergence due to the ongoing influence of $g \odot \varepsilon_t$.

Graph Analogously to SCMs, we can associate a directed graph $\mathcal{G}_{\mathcal{D}}$ to a DE \mathcal{D} , consisting of nodes V_i representing individual components $X^{(i)}$, while an edge (V_i, V_j) is present if $\exists k > 0$ such that $\partial_{X^{(i)}} X_{t+k}^{(j)} \neq 0$ in \mathcal{D} .²

2.3 Vector Autoregressive models

w

In this paper, we focus on a specific type of SDE, the VAR model (Kilian and Lütkepohl 2017).

Equation Consider a *d*-dimensional vector-valued stationary time series $\{X_0, \ldots, X_T\}$ generated by a VAR model with *lag p*, where a lag represents the number of previous time steps used to predict the current value of each variable. Specifically, the VAR(*p*) model is defined by

$$\boldsymbol{X}_t = \boldsymbol{\nu} + \boldsymbol{A}_1 \boldsymbol{X}_{t-1} + \dots + \boldsymbol{A}_p \boldsymbol{X}_{t-p} + \boldsymbol{u}_t, \quad (3)$$

where ν is a *d*-dimensional vector of intercept terms, $\{A_i\}_{i=1}^p$ are $(d \times d)$ matrices and u_t is a *d*-dimensional white noise term. If the process X_t is stable and stationary (Hamilton 1994), Equation 3 can also be written as

$$\begin{aligned} \boldsymbol{A}(L)\boldsymbol{X}_t &= \boldsymbol{\nu} + \boldsymbol{u}_t, \\ \text{ith} \quad \boldsymbol{A}(L) &:= \boldsymbol{I}_d - \boldsymbol{A}_1 L - \dots - \boldsymbol{A}_p L^p, \end{aligned} \tag{4}$$

where *L* is the *lag operator* such that $LX_t \equiv X_{t-1}$ and I_d is a *d*-dimensional identity matrix.

A key limitation of VARs is the inability to interpret the system in causal terms since the components of u_t are crosscorrelated and act as hidden confounders. A common approach to overcome this problem is to orthogonalize the noise terms. In this context, the process of *causal discovery*, i.e., inferring the causal structure of the data, is analogous to the one of SCMs (Hyvärinen et al. 2010; Moneta et al. 2013; Geiger et al. 2015; Malinsky and Spirtes 2018), and involves identifying a triangular matrix \hat{A}_0 such that $\boldsymbol{\varepsilon}_t = \hat{A}_0 \boldsymbol{u}_t$ consists of mutually uncorrelated elements. The transformed VAR, commonly known as the Structural VAR (SVAR) model in the literature (Kilian and Lütkepohl 2017), is defined by $\hat{A}_0 X_t = \hat{A}_0 \nu + \hat{A}_1 X_{t-1} + \dots + \hat{A}_p X_{t-p} + \varepsilon_t$, where $\hat{A}_i = \hat{A}_0 A_i$. From a modeling perspective, VAR and SVAR are equivalent, as any SVAR can be expressed in its reduced-form VAR by computing $A_i = \hat{A}_0^{-1} \hat{A}_i$ for i = 0, ..., p in Eq. 3. Notably, choosing one over the other does not affect its causal interpretation, provided that \hat{A}_0 is known. For simplicity, in this work, we adopt the VAR notation, to introduce a novel framework for causal inference over time, which complements the SVAR's causal discovery approach.

Graph An edge (V_i, V_j) is present iff $\exists k. A_k[i, j] \neq 0$.

¹Unlike in acyclic SCMs, $\mathbf{PA}^{(i)}$ loses its hierarchical interpretation since two nodes can be mutually parents.

²Depending on the type of DE, the derivative must be evaluated with respect to f_i or with respect to both f_i and g_i (see Equation 2).

3 Causal perspective on Discrete-time Stochastic Processes

This section provides the theoretical basis for causal inference over time. First, we formally define causal interventions on SDEs (§3.1). Then, we show how a SCM can be considered a compressed description of the asymptotic behavior of an underlying dynamical system (§3.2).

3.1 Causal Interventions on SDEs

We define an intervention \mathcal{I} on a SDE \mathcal{D} as a modification of one or more component equations denoted by the mapping:

$$\mathcal{I}: f_{i}(\mathbf{PA}_{< t}^{(i)}) + g_{i}(\mathbf{PA}_{< t}^{(i)}) \odot \varepsilon_{t} \longmapsto$$

$$\tilde{f}_{i}(\tilde{\mathbf{PA}}_{< t}^{(i)}) + \tilde{g}_{i}(\tilde{\mathbf{PA}}_{< t}^{(i)}) \odot \varepsilon_{t}, \ \forall t \ge t_{\mathcal{I}}$$
(5)

where $\tilde{\mathbf{PA}}_{<t}^{(i)} \subseteq \mathbf{PA}_{<t}^{(i)}$. Unlike SCMs, the intervention applies starting from a specific time $t_{\mathcal{I}}$. In other words, the process follows the original equations for $t < t_{\mathcal{I}}$ and the modified ones for $t \ge t_{\mathcal{I}}$. We denote the modified SDE as $\mathcal{D}^{\mathcal{I}}$ to generalize Eq. 1 to account for time. To differentiate between interventions on a SCM \mathcal{M} and on a SDE $\mathcal{D}, \mathcal{I}_{\mathcal{M}}$ and $\mathcal{I}_{\mathcal{D}}$ will be respectively adopted when necessary.

Definition 1 (Causal Effect over time (CE_t)). Let X be a solution of a specific SDE D. We define the causal effect at time t of an intervention I as

$$\mathsf{CE}_t^{\mathcal{I}} := \mathbb{E}[\boldsymbol{X}_t^{\mathcal{I}} - \boldsymbol{X}_t | \boldsymbol{X}_{< t_{\mathcal{I}}}], \tag{6}$$

where $X^{\mathcal{I}}$ is the solution of the modified SDE $\mathcal{D}^{\mathcal{I}}$.

The interpretation of $CE_t^{\mathcal{I}}$ is closely related to the causal effect of an intervention on a SCM (Eq. 1), $CE^{\mathcal{I}}$: while the latter measures the causal effect of an exogenous intervention, CE_t does so for any time step t of the DSP, i.e., *it measures the causal effect of an intervention over time*. Importantly, as we will show in the next section, $CE_t^{\mathcal{I}} \rightarrow CE^{\mathcal{I}}$ as $t \rightarrow \infty$, i.e., there is an asymptotic correspondence between the two quantities.

3.2 Mapping SDEs to SCMs

Given an SDE \mathcal{D} and its solution X, we study the conditions on \mathcal{D} such that: i) X_t converges in distribution to X_∞ as $t \to \infty$; and ii) there exists an SCM \mathcal{M} such that X_∞ is a solution of \mathcal{M} and, for every intervention \mathcal{I} , it holds that $(X_\infty)^{\mathcal{I}_{\mathcal{M}}} = (X^{\mathcal{I}_{\mathcal{D}}})_\infty$. While i) is automatically satisfied by any finite memory stationary process, ii) requires more careful analysis, as discussed below.

A negative result from (Janzing, Rubenstein, and Schölkopf 2018) Consider the stable bivariate system defined by the equations $X_t = \varepsilon_t^x, Y_t = 0.5 \cdot X_{t-1} + \varepsilon_t^y$. For every t, X_t and Y_t are independent of each other. Consequently, the joint distribution $p(X_t, Y_t)$ cannot capture the causal dependencies of the system (X causes Y). The lack of causal information in the cross-sectional dimension arises because the variables are *localized in time*; their values change rapidly, leading to minimal or no correlation with their past values. On this specific point Janzing, Rubenstein, and Schölkopf (2018) provide an explicit negative result:



Figure 1: \mathcal{T} -transformation transfers causal information from the temporal to the cross-sectional dimension, and thus to the joint distribution $P(\mathbf{Z}_t)$. The diagram commutes, i.e., red and blue paths produce the same result.

without first making the variables de-localized in time, there is no SCM that can capture the SDE. In fact, our definition of intervention (Eq. 5) *acts on a variable of the system for a prolonged and indefinite period.*

 \mathcal{T} -transformation To overcome this limitation, (Janzing, Rubenstein, and Schölkopf 2018) propose a transformation of X_t based on a frequency analysis of the time series.³ Instead, our choice is inspired by the long-run normalized mean via the transformation \mathcal{T} : DSP \mapsto DSP defined by

$$\mathcal{T}(\boldsymbol{X})_t := \boldsymbol{Z}_t = \boldsymbol{\mu} + \frac{1}{\sqrt{t}} \sum_{i=1}^t (\boldsymbol{X}_i - \boldsymbol{\mu}), \quad (7)$$

where $\mu := \mathbb{E}[X]$.⁴ Moreover, $\mathbb{E}[Z_{\infty}] = \mathbb{E}[X_{\infty}] = \mu$ so that for every intervention \mathcal{I} , $CE_{\infty}^{\mathcal{I}}$ (Eq. 6) yields the same values. However, unlike X, Z can be mapped into an SCM that precisely models *its distribution shift over any intervention*, thereby satisfying property ii), represented as the commutativity of the diagram in Fig. 1.

It is important to clarify that Z_t is not the process of interest, and the focus of the causal analysis remains on X_t . However, due to the equivalence of long-run causal effects calculated in both processes, and the ability to associate Z_t with the SCM that models these effects, Z_t serves as a convenient intermediate mathematical tool. To demonstrate how this transformation ensures these desirable properties, we will focus on the subclass of linear systems, particularly on VAR models. The reason for this choice is twofold. First, linear models, despite their simplicity, are still on par performance-wise with state-of-the-art Machine-Learning based forecast techniques (Toner and Darlow 2024), in particular when dealing with stochastic time series (Parmezan, Souza, and Batista 2019). Second, the mathematical treatment of interventions and the estimation of causal effects is particularly straightforward to implement and interpret, making this a useful first step for a possible extension to the nonlinear case.

³Our Z_t (Eq. 7) can also be interpreted as a form of discrete Fourier transform of the time-series $X_{1:t}$ with frequency zero.

⁴The expectation here is taken over time as well. Nonetheless, for stationary processes, this simplifies to $\mathbb{E}[\mathbf{X}] = \mathbb{E}[\mathbf{X}_t]$ for all t.

4 From Vector Autoregressive models to Structural Causal Models

In this section, we show that linear SCMs can model the long-term effects of stable VARs, explaining the properties of its DSP equilibrium (§4.1). Then, we provide implementations of two types of causal interventions, leveraging the strengths of the VAR framework (§4.2). Finally, we discuss the practical implications of our theoretical results (§4.3).

4.1 Mapping from VARs to SCMs

We provide the explicit mapping from VARs to linear SCMs in the following theorem (proved in App. B.2).

Theorem 1. Given a stable VAR(p) \mathcal{D} defined by Eq. 3, there exists a linear SCM \mathcal{M} with structural equations ⁵

$$\hat{X} = \hat{A}\hat{X} + \tilde{u},$$

where $\tilde{A} := [A_1 + \dots + A_p]$ and $\tilde{u} \sim \mathcal{N}(\mathbf{0}, \Sigma_{\tilde{u}}),$ (8)

such that, given the transformation $Z_t = \frac{1}{\sqrt{t}} \sum_{i=1}^t X_i$, the following properties hold:

- 1. \mathcal{D} and \mathcal{M} share the same *causal graph*, i.e., $\mathcal{G}_{\mathcal{M}} = \mathcal{G}_{\mathcal{D}}$;
- 2. The *observational distribution* induced by \mathcal{D} at equilibrium $p(\mathbf{Z}_{\infty})$ is equal to the one induced by $\mathcal{M}, p(\tilde{\mathbf{X}})$;
- 3. The *interventional distribution* $p(\mathbf{Z}_{\infty}^{\mathcal{I}})$ is equal to the one induced by the same intervention on \mathcal{M} , $p(\tilde{\mathbf{X}}^{\mathcal{I}_{\mathcal{M}}})$.

Remark. Note that due to the influence of time in VARs, the equivalent SCMs at equilibrium, while linear, may lead to cycles in the causal graph (see, e.g., Fig. 2c) and correlations between the exogenous variables, captured by the full covariance matrix $\Sigma_{\bar{u}}$ in Eq. 8. Note also that the above Theorem implies that there is a direct relationship between interventions on DSPs, \mathcal{I} , and interventions on SCMs, here denoted by $\mathcal{I}_{\mathcal{M}}$. Refer to App. B.2 for further details.

4.2 Implementation of Causal Interventions

Different application scenarios may need different types of interventions. Consider a government's fiscal policy. In such a setting, a feasible approach would be to implement an *additive intervention* in the form of an annual tax increase of, e.g., 300 euros per household on top of existing taxes. Alternatively, in other scenarios, e.g., when studying the effect of the *key European Central Bank's interest rate* (Belke and Polleit 2007), a more natural choice is to implement a *forcing intervention* that enforces the convergence of an observed variable (e.g., interest rate) to a target value. In the following, we propose an implementation for VARs of these two forms of interventions, showing their effects on the system and discussing their stability conditions.

Additive Interventions Given a stable VAR(p) as in Eq. 4, we define an additive intervention \mathcal{I}_a at time $t_{\mathcal{I}}$ with force F as the mapping:

$$\mathcal{I}_a: \mathbf{A}(L)\mathbf{X}_t = \mathbf{\nu} + \mathbf{u}_t \quad \longmapsto \\ \mathbf{A}(L)\mathbf{X}_t = \mathbf{\nu} + \mathbf{u}_t + \mathbb{I}(t \ge t_{\mathcal{I}})\mathbf{F}, \tag{9}$$

where $\mathbb{I}(t \ge t_{\mathcal{I}})$ is the indicator function, which equals 1 if $t \ge t_{\mathcal{I}}$, otherwise 0. In other words, we perform a translation while keeping the process dynamics unchanged. In such case, the temporal causal effect CE_t is deterministic and takes values $CE_t = 0$ for $t < t_{\mathcal{I}}$ while, for $k \ge 0$:

$$\mathtt{CE}_{t_{\mathcal{I}}+k} = \sum_{l=0}^{k} \Phi_l \boldsymbol{F},$$

where Φ is the *impulse response function* of the VAR model. Refer to App. B.1 for further details.

Remark. For this type of intervention, CE_t is deterministic and does not depend on the specific trajectory. The same property can be observed on the linear SCM associated with the process, defining the intervention in a similar way: $\tilde{X} = \tilde{A}\tilde{X} + \tilde{u}$ changes into $\tilde{X} = \tilde{A}\tilde{X} + F + \tilde{u}$.

Stability Additive intervention preserves the stability regardless of the value of F, since A(L) does not change. See App. B.1 for stability conditions of VARs.

Forcing Interventions We define a forcing intervention \mathcal{I}_f at time $t_{\mathcal{I}}$ with force F and target value \hat{X} as:

$$\mathcal{I}_{f}: \boldsymbol{A}(L)\boldsymbol{X}_{t} = \boldsymbol{\nu} + \boldsymbol{u}_{t} \quad \longmapsto \\ \boldsymbol{A}(L)\boldsymbol{X}_{t} = \boldsymbol{\nu} + \boldsymbol{u}_{t} + \mathbb{I}(t \ge t_{\mathcal{I}})\boldsymbol{F} \odot (\hat{\boldsymbol{X}} - \boldsymbol{X}_{t}).$$
(10)

We assume F to be positive in each component. This intervention acts as an attraction towards \hat{X} , and F modulates the intensity of the attraction force. Applying an intervention on a single component $X^{(i)}$ toward the fixed value \hat{X} and letting $F^{(i)} \to +\infty$ yields the do operator $do(X^{(i)} = \hat{X})$. We refer to (Mooij, Janzing, and Schölkopf 2013) for a detailed discussion of this point.

Stability Forcing interventions \mathcal{I}_f perturb the system dynamics by modifying the operator A(L). Specifically, by shifting the term $F \odot X_t$ to the left of the equation and rewriting it in matrix form as $F_{diag}X_t$, we obtain $\tilde{A}(L) := A(L) + F_{diag}$. Hence, the stability of the intervened system is not guaranteed (we provide an example in App. B.3), and it is necessary to verify *that all the eigenvalues of* $\tilde{A}(L)$ are still inside the unit circle. Intuitively, the stability of an observational system often relies on negative feedback loops. Fixing one variable can disrupt this balance, leading to runaway behavior. For example, turning off a pressure release valve in a pressurized tank can cause the pressure to build up uncontrollably, eventually leading to an explosion.

4.3 Practical implications

Causal queries Our formulation of causal interventions on VARs differs from the standard approach based on Granger causality by being closer to that of SCMs. Consequently, it enables the generalization of interventional and counterfactual queries to account for time (see App. C). That is, it allows for answering the following causal questions:

• Forecasted Interventions What are the expected effects on an individual trajectory (or a population) when intervening in the present, and how do they vary over time?

⁵For simplicity we set $\nu = 0$, i.e., we assume $\mathbb{E}[\mathbf{X}_t] = 0$. The theorem applies in the general case up to a translation of both the VAR and the associated SCM.

• **Retrospective Counterfactuals** What would have happened to an individual trajectory if an intervention had been applied at a specific point in the past? What state would it be in now?

Both causal queries acquire a meaning embedded in the temporal dimension in terms of *forecasting for the future* (§5.2) and *retrospection for the past* (App. D.3), respectively.

Expressiveness and universality VARs, despite their linearity, possess a high level of expressiveness (Kilian and Lütkepohl 2017). In fact, the Wold decomposition Theorem (Wold 1938) implies that the dynamics of *any purely nondeterministic covariance-stationary process can be approximated arbitrarily well by an autoregressive model, making them universal approximators*. In practice, linear autoregressive models are broadly used in time-series analysis. Yet, we intend to explore non-linear DSPs in future work, as they may lead to better convergence rates and allow for causal *interpretation of a broader family of dynamical models*.

Feedback loops To properly understand complex systems, it is often useful to model feedback loops between their variables. Time-series models naturally capture this property, while SCMs require significant reformulation. The theory of cyclic SCMs has seen a significant advancement in recent years (Bongers et al. 2021), but practical approaches, both for causal discovery and causal inference, are still underdeveloped (Bongers et al. 2016; Lorbeer and Mohsen 2023). Our formalization of causal inference on VARs is a step forward in this direction.

Fitting VARs estimation is typically performed using ordinary least squares. Various alternative methods are available, both in terms of constrained optimization (e.g., to use prior knowledge about some coefficients of the VAR matrices (Sims 1980)) and within a Bayesian framework (Koop, Korobilis et al. 2010). Refer to (Lütkepohl 2005, chapters 3,4,5) for a comprehensive discussion. Importantly, although VARs are most commonly used on time-series data (i.e., data from one single unit across a period of time), there are approaches tailored to the analysis of panel data (i.e., the evolution of many units over time) (Sigmund and Ferstl 2021); and cross-sectional data (i.e., many individuals at a single point of time), provided that they have at least some proxy variables of time (Deaton 1985). Such approaches open up a promising line of future work that can further generalize VARs applicability for causal reasoning over time.

5 Empirical evaluation

In this section, we evaluate VAR models' accuracy and expressiveness in multivariate time series, focusing on two forecasting dimensions: observational (§5.1) and interventional (§5.2). Additional results and in-depth descriptions can be found in App. D.

Datasets We rely on two synthetic datasets, German⁶ and Pendulum, and the real-world Census dataset⁷. German sim-

⁷https://data.census.gov/



Figure 2: **Causal graphs**. The causal graph for (a) and (b) is known, while for (c), it is assumed. In (b), nodes are labeled with the initials of each feature: Expertise, Responsibility, Loan Amount, Duration, Income, Savings, and Credit Score. In (c), 0 - 14, 15 - 64, and 65 - 99 represent age groups.

ulates a loan approval scenario with seven variables. Pendulum is a two-variable system where $X^{(1)}$ operates as a stabilizer for $X^{(2)}$, which exhibits a divergent dynamic. Census includes demographic variables across three age groups, along with migration, birth, and death rates from 1992 to 2023 for 50 countries. Fig. 2 illustrates the causal graphs for all datasets. See App. D for further details.

Metrics We measure the discrepancy between the *h*-steps forecast $\hat{X}_{t+h} | X_{<t}$ and the true value X_{t+h} on the test set X_{test} . We report Mean Absolute Error (MAE) focusing on the target variables (i.e., *Credit Score* for German, $X^{(1)}$ for Pendulum, and age groups for Census). See App. D.1 for other metrics. All results are averaged over ten runs.

5.1 Observational Forecasting

Baselines We compare VAR with three relevant works: i) DLinear (Zeng et al. 2023), a decomposition-based linear model that separates trend and seasonal components; ii) TSMixer (Chen et al. 2023), a Multi-layer Perceptron (MLP) based model that focuses on mixing time and feature dimensions; iii) TiDE (Das et al. 2023), a MLP based encoderdecoder model. To assess the effectiveness of the forecasting methods, we introduce an observational oracle forecaster that has full knowledge about the true data generating process and produces the optimal predictor, i.e., $\hat{X}_{t+h}|X_{< t} = \mathbb{E}[X_{t+h}|X_{< t}]$.

How does the VAR performance compare with SOTA models for forecasting multivariate time series? The observational forecasting results in Table 1 show performance across varying data sizes (i.e., number of instances) and forecast horizons for all datasets. VAR emerges as the top-performing model, consistently matching or closely approaching Oracle's scores for all datasets. DLinear usually achieves predictive accuracy close to VAR for 1-step forecasts, presumably due to the common linear nature of both models. TiDE and TSMixer consistently underperform compared to VAR and DLinear for German and Pendulum. For all models (including Oracle), performance on the Pendulum dataset is uniformly worse than on the German, highlighting the greater challenge in forecasting given the system's stronger stochasticity and variables changing more rapidly over time. On Census, VAR and TiDE provide the best results, TiDE slightly outperforming VAR in 5-step horizon.

⁶This dataset is inspired on https://archive.ics.uci.edu/dataset/ 144/statlog+german+credit+data



Figure 3: Additive Intervention. (Left) Intevention on *Expertise* with F = 0.2. (Right) Effect on *Credit Score*. Shaded regions in both plots denote 95% confidence bounds.

Table 1: **Observational Forecasting.** MAE scores (*lower is better*) for VAR, DLinear (Zeng et al. 2023), TiDE (Das et al. 2023) and TSMixer (Chen et al. 2023). Results averaged over 10 runs. Due to space limitations, standard deviations are reported in App. D.1. Best model in bold, Oracle in typewriter. For Census, size equals the number of countries times the number of years.

			Observational Forecasting				
Dataset	Size	Horizon	Oracle	VAR	DLinear	TiDE	TSMixer
German Pendulum	100	1	.004 .042	.008 .043	.009 . 043	.011 .218	.014 .217
German Pendulum		10	.014 .399	.055 .420	. 055 .440	$\begin{array}{c} .094 \\ 1.43 \end{array}$	$.139 \\ 1.43$
German Pendulum	500	1	.004 .042	.004 .042	.004 .042	.011 .218	.014 .217
German Pendulum		10	.014 .399	.015 .401	. 015 .405	.093 1.43	$.135 \\ 1.43$
Census	50×32	1 5	-	. 001 .017	.006 .025	.001 .014	.008 .024

5.2 Interventional Forecasting

We evaluate the causal VAR's forecast in estimating the causal effects on German. See App. D.2 for other datasets.

Baselines Since state-of-the-art methods do not allow computing the causal effect of interventions on dynamical systems, we use an oracle forecaster as a benchmark for theoretically optimal performance. Specifically, the ground truth values are estimated as $CE_{t+h} = \mathbb{E}[\mathbf{X}_{t+h}^{\mathcal{I}} - \mathbf{X}_{t+h} | \mathbf{X}_{< t}]$, while the predicted values from the proposed VAR framework as $\hat{CE}_{t+h} = (\hat{\mathbf{X}}_{t}^{\mathcal{I}} - \hat{\mathbf{X}}_{t+h}) | \mathbf{X}_{< t}$.

Interventions We perform causal interventions on the root node *Expertise* and observe the effect on the target variable *Credit Score*. For the additive case, we apply F = 0.2, while for forcing, we use F = 1 with a target value of $\hat{E} = 5$.



Figure 4: Forcing Intervention. (Left) Intervention on *Expertise* with F = 1 and target $\hat{E} = 5$. (Right) Effect on *Credit Score*. Note that confidence bounds are narrower w.r.t. Fig. 3.

Table 2: **Interventional Forecasting.** MAE scores (*lower is better*) for the proposed causal VAR framework on the German dataset. Results averaged over 10 runs, with standard deviation in subscript. Scores are scaled by a factor of 10^2 to ease readability.

			Interventional Forecasting		
Dataset	Size	Horizon	Additive	Forcing	
German	100	1 10	.000 _{.000} .043 _{.028}	$.000_{.000}$ $.364_{.297}$	
	500	1 10	.000 _{.000} .018 _{.014}	$.000_{.000}$ $.115_{.081}$	

These values are selected for illustrative purposes such that the long-term expected value of *Expertise* is the same for both interventions (i.e., 5). See App. D.2 for other variants.

How do additive and forcing interventions affect the system dynamics? Fig. 3 illustrates the additive intervention. *Expertise* is a variable that typically necessitates several years for acquisition in practical scenarios. The causal VAR accurately captures such delayed impact as its effect on *Credit Score* appears after several years. As the system maintains its dynamic characteristics unchanged, the forecasted covariance remains the same even after the intervention. Fig. 4 shows the forcing intervention, where the interventional forecasting exactly aligns with the target value. Moreover, we stress that even for a low value of F, the forcing intervention resembles a do-intervention (shrinking the variances significantly) even though theoretically convergence is guaranteed only for $F \to \infty$.

How accurate is the causal VAR framework in estimating the causal effect of interventions over time? Table 2 summarizes results on interventional forecasting, showing errors with varying data size and forecast horizons. At 1step, both interventions lead to perfect performance since *Credit Score* is a slow-changing variable and requires at least 3 time steps for an intervention to take effect. At 10-step, our causal effect estimates remain highly accurate.



Figure 5: **German**. Effect of increasing *Expertise* on *Credit Score*. (Left) The time each loan applicant takes to cross or not the acceptance threshold. The histogram shows the distribution of crossing times. (Right) Comparison of two loan applicants, i.e., trajectories, with similar scores at intervention time. After the intervention, they diverge significantly, with only an applicant being accepted at the maximum time. Forecasts are dashed for observational and dotted for interventional.

6 Use cases

In this section, we show real-world scenarios where estimating causal effects over time represents a useful step toward a realistic modeling of the phenomena. We focus on the German and Census datasets, presenting two analyses for each. For German, Fig. 5 illustrates the effect of increasing *Expertise* by F = 0.38 on *Credit Score*.

German 1 – Same intervention, different outcomes. The left panel of Fig. 5 shows the distribution of the time required for loan applicants to cross or not the acceptance threshold after the intervention. Access to this information allows quantification of intervention efficacy, identification of credit-building patterns, and infer the key factors influencing loan eligibility outcomes. It can also inform the recommendation of actions (e.g. in algorithmic recourse (Karimi et al. 2022)) within a reasonable timeframe, fostering trust in the system and promoting user acceptance.

German 2 – Similar cross-sectional values, different causal effects over time. The right panel of Fig. 5 shows trajectories that, while seeming similar at a given time, may have significantly divergent historical and future behaviors. For instance, the green trajectory may have autonomously crossed the threshold without intervention, whereas in the case of the yellow one, the applied intervention may be inadequate to ensure the desired outcome. Such divergence highlights the importance of moving beyond models that rely only on cross-sectional data, motivating the need for techniques, such as the proposed causal VAR framework, that capture individual applicant behavior over time.

For Census, Fig. 6 presents two additive interventions across all countries.

Census 1 – Impact of Births on Avg. Age. The left panel of Fig. 6 shows the increase on *Births* with F = 0.004 and its effect on the population's average age (computed as a weighted mean of age groups). The force value means that *Births* increase by 0.4% of each country's total population



Figure 6: **Census**. Additive intervention across all countries. (Left) Intervention on *Births* with F = 0.004 and its effect on the population's average age. (Right) Intervention on *Migration* with F = 0.04 and its effect on *Births*. In both plots, we highlight two countries (one above and one below the mean) to illustrate the differences after the intervention. Forecasts are dashed for observational and dotted for interventional.

every year. Examining how they affect population age over time could allow policymakers to identify which countries might benefit most from specific types of demographic interventions and evaluate the long-term viability of systems. We can also observe that the intervention in Japan causes a more evident decrease in the average age than in Chile.

Census 2 – Impact of Migration on Births. The right panel of Fig. 6 reports how a 4% growth in *Migration* w.r.t. the total population influences *Birth* rates. We observe that increased migration leads to a rise in births. However, its impact is less evident (observational and interventional forecasting trajectories are closer) compared to the result on the average age shown in Fig. 6 (left panel).

7 Concluding remarks

In this work, we have established a link between discretetime dynamical systems at equilibrium and SCMs. Moreover, we have provided an explicit procedure for mapping VARs to linear SCMs and demonstrated that, under specific model stability conditions, interventions on the dynamical system and the SCM yield equivalent results. To conduct causal inference over time, we have introduced two classes of interventions (additive and forcing) for VARs.

Limitations When systems exhibit strongly nonlinear dynamics, linear VARs may prove less effective than alternative nonlinear approaches. Moreover, our framework requires prior knowledge of the causal graph. In scenarios where this information is lacking, the process of causal discovery can present significant challenges.

Future work We will investigate the use of non-linear DSPs, as they may lead to better convergence rates and allow for causal interpretation of a broader family of dynamical models. Moreover, our work opens several interesting research directions (§4.3 for concrete examples) and applications (e.g. causal inference over time in high-dimensional contexts such as climate science).

References

Anderson, T. 1994. The Statistical Analysis of Time Series. *Wiley Series in Probability and Statistics*.

Arnold, L.; Jones, C. K.; Mischaikow, K.; Raugel, G.; and Arnold, L. 1995. *Random dynamical systems*. Springer.

Bareinboim, E.; Correa, J. D.; Ibeling, D.; and Icard, T. 2022. On Pearl's Hierarchy and the Foundations of Causal Inference. In *Probabilistic and Causal Inference*, volume 36 of *ACM Books*, 507–556. ACM.

Belke, A.; and Polleit, T. 2007. How the ECB and the US Fed set interest rates. *Applied Economics*, 39(17): 2197–2209.

Bhattacharya, R.; and Majumdar, M. 2003. Random dynamical systems: a review. *Economic Theory*, 23: 13–38.

Bongers, S.; Blom, T.; and Mooij, J. M. 2018. Causal modeling of dynamical systems. arXiv preprint arXiv:1803.08784.

Bongers, S.; Forré, P.; Peters, J.; and Mooij, J. M. 2021. Foundations of structural causal models with cycles and latent variables. *The Annals of Statistics*, 49(5): 2885–2915.

Bongers, S.; Peters, J.; Schölkopf, B.; and Mooij, J. M. 2016. Structural Causal Models: Cycles, Marginalizations, Exogenous Reparametrizations and Reductions. *CoRR*, abs/1611.06221.

Chen, S.; Li, C.; Arik, S. Ö.; Yoder, N. C.; and Pfister, T. 2023. TSMixer: An All-MLP Architecture for Time Series Forecast-ing. *Trans. Mach. Learn. Res.*, 2023.

Das, A.; Kong, W.; Leach, A.; Mathur, S.; Sen, R.; and Yu, R. 2023. Long-term Forecasting with TiDE: Time-series Dense Encoder. *Trans. Mach. Learn. Res.*, 2023.

Deaton, A. 1985. Panel data from time series of cross-sections. *Journal of econometrics*, 30(1-2): 109–126.

Edwards, R. M.; Lee, K. Y.; and Ray, A. 1992. Robust optimal control of nuclear reactors and power plants. *Nuclear Technology*, 98(2): 137–148.

Geiger, P.; Zhang, K.; Schoelkopf, B.; Gong, M.; and Janzing, D. 2015. Causal Inference by Identification of Vector Autoregressive Processes with Hidden Components. In Bach, F.; and Blei, D., eds., *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, 1917–1925. Lille, France: PMLR.

Granger, C. W. 1969. Investigating causal relations by econometric models and cross-spectral methods. *Econometrica: journal of the Econometric Society*, 424–438.

Hamilton, J. D. 1994. Time Series Analysis.

Herzen, J.; Lässig, F.; Piazzetta, S. G.; Neuer, T.; Tafti, L.; Raille, G.; Van Pottelbergh, T.; Pasieka, M.; Skrodzki, A.; Huguenin, N.; et al. 2022. Darts: User-friendly modern machine learning for time series. *Journal of Machine Learning Research*, 23(124): 1–6.

Hildebrand, F. B. 1987. *Introduction to numerical analysis*. Courier Corporation.

Hyvärinen, A.; Zhang, K.; Shimizu, S.; and Hoyer, P. O. 2010. Estimation of a Structural Vector Autoregression Model Using Non-Gaussianity. *J. Mach. Learn. Res.*, 11: 1709–1731.

Janzing, D.; Rubenstein, P.; and Schölkopf, B. 2018. Structural causal models for macro-variables in time-series. *arXiv preprint arXiv:1804.03911*.

Jornet, M. 2023. Theory and methods for random differential equations: a survey. *SeMA Journal*, 80(4): 549–579.

Karimi, A.-H.; Barthe, G.; Schölkopf, B.; and Valera, I. 2022. A survey of algorithmic recourse: contrastive explanations and consequential recommendations. *ACM Computing Surveys*, 55(5): 1–29.

Kilian, L.; and Lütkepohl, H. 2017. *Structural vector autoregressive analysis*. Cambridge University Press.

Kingma, D. P.; and Ba, J. 2015. Adam: A Method for Stochastic Optimization. In *ICLR (Poster)*.

Koop, G.; Korobilis, D.; et al. 2010. Bayesian multivariate time series methods for empirical macroeconomics. *Foundations and Trends*® *in Econometrics*, 3(4): 267–358.

Lasota, A.; and Mackey, M. C. 1989. Stochastic perturbation of dynamical systems: The weak convergence of measures. *Journal of Mathematical Analysis and Applications*, 138(1): 232–248.

Lim, B.; and Zohren, S. 2021. Time-series forecasting with deep learning: a survey. *Philosophical Transactions of the Royal Society A*, 379(2194): 20200209.

Lorbeer, B.; and Mohsen, M. 2023. Comparative Study of Causal Discovery Methods for Cyclic Models with Hidden Confounders. In *CogMI*, 103–111. IEEE.

Łoskot, K.; and Rudnicki, R. 1995. Limit theorems for stochastically perturbed dynamical systems. *Journal of applied probability*, 32(2): 459–469.

Lütkepohl, H. 2005. *New introduction to multiple time series analysis*. Springer Science & Business Media.

Malinsky, D.; and Spirtes, P. 2018. Causal Structure Learning from Multivariate Time Series in Settings with Unmeasured Confounding. In *CD@KDD*, volume 92 of *Proceedings of Machine Learning Research*, 23–47. PMLR.

Moneta, A.; Chlaß, N.; Entner, D.; and Hoyer, P. 2011. Causal search in structural vector autoregressive models. In *NIPS Mini-Symposium on Causality in Time Series*, 95–114. PMLR.

Moneta, A.; Entner, D.; Hoyer, P. O.; and Coad, A. 2013. Causal inference by independent component analysis: Theory and applications. *Oxford Bulletin of Economics and Statistics*, 75(5): 705–730.

Montgomery, D. C.; Jennings, C. L.; and Kulahci, M. 2015. *Introduction to time series analysis and forecasting*. John Wiley & Sons.

Mooij, J. M.; Janzing, D.; and Schölkopf, B. 2013. From Ordinary Differential Equations to Structural Causal Models: the deterministic case. In Nicholson, A. E.; and Smyth, P., eds., *Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence, UAI 2013, Bellevue, WA*, USA, August 11-15, 2013. AUAI Press.

Parmezan, A. R. S.; Souza, V. M.; and Batista, G. E. 2019. Evaluation of statistical and machine learning models for time series prediction: Identifying the state-of-the-art and the best conditions for the use of each model. *Information sciences*, 484: 302–337.

Pearl, J. 2009. Causality. Cambridge university press.

Runge, J.; Bathiany, S.; Bollt, E.; Camps-Valls, G.; Coumou, D.; Deyle, E.; Glymour, C.; Kretschmer, M.; Mahecha, M. D.; Muñoz-Marí, J.; et al. 2019. Inferring causation from time series in Earth system sciences. *Nature communica-tions*, 10(1): 2553.

Sigmund, M.; and Ferstl, R. 2021. Panel vector autoregression in R with the package panelvar. *The Quarterly Review of Economics and Finance*, 80: 693–720.

Sims, C. A. 1980. Macroeconomics and reality. *Econometrica: journal of the Econometric Society*, 1–48.

Toner, W.; and Darlow, L. N. 2024. An Analysis of Linear Time Series Forecasting Models. In *Forty-first International Conference on Machine Learning*.

Wold, H. 1938. A study in the analysis of stationary time series. Ph.D. thesis, Almqvist & Wiksell.

Wunsch, G.; Russo, F.; Mouchart, M.; and Orsi, R. 2022. Time and causality in the social sciences. *Time & Society*, 31(2): 177–204.

Zeng, A.; Chen, M.; Zhang, L.; and Xu, Q. 2023. Are Transformers Effective for Time Series Forecasting? *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(9): 11121–11128.

Reproducibility Checklist

This paper:

- Includes a conceptual outline and/or pseudocode description of AI methods introduced: **yes** (§4)
- Clearly delineates statements that are opinions, hypotheses, and speculation from objective facts and results: **yes**
- Provides well-marked pedagogical references for lessfamiliar readers to gain background necessary to replicate the paper: **yes**

Does this paper make theoretical contributions? yes

- All assumptions and restrictions are stated clearly and formally. **yes**
- All novel claims are stated formally (e.g., in theorem statements). **yes**
- Proofs of all novel claims are included. yes
- Proof sketches or intuitions are given for complex and/or novel results. **yes**
- Appropriate citations to theoretical tools used are given. **yes**
- All theoretical claims are demonstrated empirically to hold. **yes**
- All experimental code used to eliminate or disprove claims is included. **yes**

Does this paper rely on one or more datasets? yes

• A motivation is given for why the experiments are conducted on the selected datasets **yes**

- All novel datasets introduced in this paper are included in a data appendix. **yes**
- All novel datasets introduced in this paper will be made publicly available upon publication of the paper with a license that allows free usage for research purposes. **yes**
- All datasets drawn from the existing literature (potentially including authors' own previously published work) are accompanied by appropriate citations. **yes**
- All datasets drawn from the existing literature (potentially including authors' own previously published work) are publicly available. **yes**
- All datasets that are not publicly available are described in detail, with an explanation why publicly available alternatives are not scientifically satisfactory. **NA**

Does this paper include computational experiments? yes

- Any code required for pre-processing data is included in the appendix. **yes**
- All source code required for conducting and analyzing the experiments is included in a code appendix. **yes**
- All source code required for conducting and analyzing the experiments will be made publicly available upon publication of the paper with a license that allows free usage for research purposes. **yes**
- All source code implementing new methods has comments detailing the implementation, with references to the paper where each step comes from **yes**
- If an algorithm depends on randomness, then the method used for setting seeds is described in a way sufficient to allow replication of results. **yes**
- This paper specifies the computing infrastructure used for running experiments (hardware and software), including GPU/CPU models; amount of memory; operating system; names and versions of relevant software libraries and frameworks. **yes**
- This paper formally describes evaluation metrics used and explains the motivation for choosing these metrics. **yes**
- This paper states the number of algorithm runs used to compute each reported result. **yes**
- Analysis of experiments goes beyond single-dimensional summaries of performance (e.g., average; median) to include measures of variation, confidence, or other distributional information. **yes**
- The significance of any improvement or decrease in performance is judged using appropriate statistical tests (e.g., Wilcoxon signed-rank). **no**
- This paper lists all final (hyper-)parameters used for each model/algorithm in the paper's experiments. **yes**
- This paper states the number and range of values tried per (hyper-)parameter during development of the paper, along with the criterion used for selecting the final parameter setting. **yes**