

Causal Inference on Time Series using Structural Equation Models

Jonas Peters*[†]
peters@stat.math.ethz.ch

Dominik Janzing*
janzing@tuebingen.mpg.de

Bernhard Schölkopf*
bs@tuebingen.mpg.de

July 24, 2012

Abstract

Causal inference uses observations to infer the causal structure of the data generating system. We study a class of functional models that we call Time Series Models with Independent Noise (TiMINo). These models require independent residual time series, whereas traditional methods like Granger causality exploit the variance of residuals. There are two main contributions: (1) *Theoretical*: By restricting the model class (e.g. to additive noise) we can provide a more general identifiability result than existing ones. This result incorporates lagged and instantaneous effects that can be nonlinear and do not need to be faithful, and non-instantaneous feedbacks between the time series. (2) *Practical*: If there are no feedback loops between time series, we propose an algorithm based on non-linear independence tests of time series. When the data are causally insufficient, or the data generating process does not satisfy the model assumptions, this algorithm may still give partial results, but mostly avoids incorrect answers. An extension to (non-instantaneous) feedbacks is possible, but not discussed. It outperforms existing methods on artificial and real data. Code can be provided upon request.

1 Introduction

We consider finitely many time series $X_t^i, i \in V$, with a maximal time order of p , that is we assume no influence from X_{t-k}^i on X_t^j for $k > p$. We further assume stationarity: the influence from X_{t-k}^i on X_t^j is required to be the same for all t . The question whether X^i is causing X^j now reads as whether there is a causal influence from some X_{t-k}^i on X_t^j , for $0 \leq k < p$. All models assume homoscedastic noise.

We first review causal inference on iid data, that is in the case with no time structure, in Section 2. Note that iid methods cannot be applied directly on time series data because a common history might introduce complicated dependencies between contemporaneous data X_t and Y_t . Motivated by the iid case, Chu and Glymour [2008] and Hyvärinen et al. [2008] propose approaches for the time series setting that include linear instantaneous effects. We describe these methods together with Granger causality in Section 3. All of them encounter similar problems: none of them are general enough to include nonlinear instantaneous effects or hidden common causes. Furthermore, when the model assumptions are violated the methods give incorrect results and one draws false causal conclusions without noticing. We propose to use time series models with independent noise (*TiMINo*) that include nonlinear and instantaneous effects. The model is based on Functional Models (also known as Structural Equation Models) and assumes X_t to be a function of all direct causes and some noise variable, the collection of which is supposed to be jointly independent. This constitutes a relatively straight-forward extension on iid methods, but we regard the benefits in the setting of time series as substantial: In Section 4 we prove that for TiMINo models the full causal structure can be recovered from the distribution. Section 5 introduces an algorithm (*TiMINo causality*) that recovers the model structure from a finite sample. It covers a broader class of models than

*Max Planck Institute for Intelligent Systems, Tübingen, Germany

[†]Seminar for Statistics, ETH Zurich, Switzerland

existing methods and can be run with any provided algorithm for fitting time series. If the data do not satisfy the assumptions, TiMINo causality remains mostly (see Section 5.3) undecided instead of drawing wrong causal conclusions. The methods are applied to simulated and real data sets in Section 6.

2 Causal inference on iid data

Inferring causal relations from observational data is challenging when interventions are not applicable. Given iid samples from $\mathbf{P}^{(X^i), i \in V}$, we try to find the underlying causal structure of the variables $X^i, i \in V$.

2.1 Directed acyclic graphs and constraint-based methods

Let $X^i, i \in V$ be a set of random variables and \mathcal{G} a directed acyclic graph (DAG) on V . The joint distribution is said to be *Markov* with respect to the DAG \mathcal{G} if each variable is independent of its non-descendants given its parents. The distribution is *faithful* with respect to \mathcal{G} if all conditional independences are entailed by the Markov assumption. Constraint-based methods [e.g. Spirtes et al., 2000] assume that the joint distribution is Markov, and faithful with respect to the true causal DAG. They show how to exploit conditional independences for reconstructing the graph \mathcal{G} , e.g. using the PC algorithm; but the graph can only be recovered up to *Markov equivalence classes*. E.g., $X \rightarrow Y$ and $Y \rightarrow X$ cannot be distinguished.

2.2 Functional models and additive noise

Functional models [Pearl, 2009] provide a different approach to the problem described above: We say $\mathbf{P}^{(X^i), i \in V}$ satisfies a *functional model* if for all $i \in V$ there exists a set of nodes $\mathbf{PA}^i \subseteq X^{V \setminus \{i\}}$, a function f_i and a noise variable N^i , such that we can write $X^i = f_i(\mathbf{PA}^i, N^i)$. (For any subset $\mathbf{A} \subset V$ we define $X_{\mathbf{A}} := \{X^j \mid j \in \mathbf{A}\}$. Additionally, we require $(N^i)_{i \in V}$ to be jointly independent and the graph obtained by drawing arrows from all elements of \mathbf{PA}^i to X^i (for each $i \in V$) to be acyclic. By restricting the function class one can identify the bivariate case: Shimizu et al. [2006] show that if $\mathbf{P}^{(X, Y)}$ allows for $Y = a \cdot X + N_Y$ with $N_Y \perp\!\!\!\perp X$ then $\mathbf{P}^{(X, Y)}$ only allows for $X = b \cdot Y + N_X$ if (X, N_Y) are jointly Gaussian ($\perp\!\!\!\perp$ stands for statistical independence). This idea has led to the extensions of nonlinear additive functions $f(x, n) = g(x) + n$ [Hoyer et al., 2009], post-nonlinear additive functions $f(x, n) = h(g(x) + n)$ [Zhang and Hyvarinen, 2009] and discrete functions [Peters et al., 2011a]. Peters et al. [2011b] show that identifiability in the bivariate case is enough for multiple variables. Mooij et al. [2009] provides practical ANM-based methods for more than two variables. Sections 4 and 5 apply these principles to time series.

3 Causal inference on time series: existing methods

For each i from a finite V , let $(X_t^i)_{t \in \mathbb{N}}$ be a time series. \mathbf{X}_t denotes the vector of time series values at time t . We call the infinite graph that contains each variable X_t^i as a node the *full time graph*. The *summary time graph* contains all $\#V$ components of the time series as vertices and an arrow between X^i and $X^j, i \neq j$, if there is an arrow from X_{t-k}^i to X_t^j in the full time graph for some k . This work addresses the following

Problem: *Given a sample $(\mathbf{X}_1, \dots, \mathbf{X}_T)$ of a multivariate time series, recover the true causal summary time graph.*

3.1 Granger causality

Granger causality [Granger, 1969] (*G-causality* for the remainder of the article) does not require complicated statistics, it is easy to implement, and it is based on the following idea: X^i does not Granger cause X^j if including the past of X^i does not help in predicting X_t^j given the past of all other time series $X^k, k \neq i$. In principle, “all other” means all other information in the world. In practice, one is limited to $X^k, k \in V$. In order to translate the phrase “does not help” into the mathematical language we need

to assume a multivariate time series model. If the data follow the assumed model, e.g. the VAR model below, G-causality is sometimes interpreted as testing whether $X_{t-h}^i, h > 0$ is independent of X_t^j given $X_{t-h}^k, k \in V \setminus \{i\}, h > 0$ [see Florens and Mouchart, 1982, Eichler, 2011, Chu and Glymour, 2008, and Section 3.2].

3.1.1 Linear Granger causality

Linear G-causality considers a VAR model: $\mathbf{X}_t = \sum_{\tau=1}^p \mathbf{A}(\tau)\mathbf{X}_{t-\tau} + \mathbf{N}_t$, where \mathbf{X}_t and \mathbf{N}_t are vectors and $\mathbf{A}(\tau)$ are matrices. For checking whether X^i G-causes X^j one fits a full VAR model M_{full} to \mathbf{X}_t and a VAR model M_{restr} to \mathbf{X}_t with the constraints $A_{\cdot i}(\tau) = 0$ for all $1 \leq \tau \leq p$ that predicts X_t^i without using X^j . Then one checks whether the reduction of the residual sum of squares (RSS) of X_t^i is significant by using the following test statistic: $T := \frac{(RSS_{\text{restr}} - RSS_{\text{full}})/(p_{\text{full}} - p_{\text{restr}})}{RSS_{\text{full}}/(N - p_{\text{full}})}$, where p_{full} and p_{restr} are the number of parameters in the respective models. For the significance test we use $T \sim F_{p_{\text{full}} - p_{\text{restr}}, N - p_{\text{full}}}$.

3.1.2 Nonlinear Granger causality

G-causality has been extended to nonlinear relationships, [e.g. Chen et al., 2004, Ancona et al., 2004]. In this paper we focus on an extension for the bivariate case proposed by Bell et al. [1996]. It is based on generalized additive models (gams) [Hastie and Tibshirani, 1990]: $X_t^i = \sum_{\tau=1}^p \sum_{j=1}^n f_{i,j,\tau}(X_{t-\tau}^j) + N_t^i$, where N_t is a $\#V$ dimensional noise vector. In order to test whether X^2 G-causes X^1 , for example with order 1, two models are fit: $X_t^1 = g_1(X_{t-1}^1) + N_t$ and $X_t^1 = g_2(X_{t-1}^1) + g_3(X_{t-1}^2) + M_t$. Bell et al. [1996] utilize the same F statistic as above; this time p_{full} and p_{restr} are the estimated degrees of freedom of the corresponding models. They refer to simulation studies by Hastie and Tibshirani [1990].

3.2 ANLTSM

Following Bell et al. [1996], Chu and Glymour [2008] introduce additive nonlinear time series models (ANLTSM for short) for performing relaxed conditional independence tests: If including one variable, e.g. X_{t-1}^1 , into a model for X_t^2 that already includes X_{t-2}^2, X_{t-1}^2 , and X_{t-2}^1 does not improve the predictability of X_t^2 , then X_{t-1}^1 is said to be independent of X_t^2 given $X_{t-2}^2, X_{t-1}^2, X_{t-2}^1$ (if the maximal time lag is 2). Chu and Glymour [2008] propose a method based on constraint-based methods like FCI [Spirtes et al., 2000] in order to infer the causal structure exploiting those conditional independence statements. The instantaneous effects are assumed to be linear and the confounders linear and instantaneous. Unfortunately, we did not find code for this method.

3.3 TS-LiNGAM

LiNGAM [Shimizu et al., 2006] infers causal graphs for linear, non-Gaussian data. It has been extended to time series by Hyvärinen et al. [2008] (for short: TS-LiNGAM). It allows for instantaneous effects, all relationships are assumed to be linear. Hidden confounders and nonlinearities may lead to wrong results.

3.4 Limitations of existing methods

The approaches described above suffer from the following methodological problems: (1) *Instantaneous effects*: The formulation of G-causality has the intrinsic problem that it cannot deal with instantaneous effects. E.g., when X_t is causing Y_t , including any of the two time series helps for predicting the other. Thus G-causality infers $X \rightarrow Y$ and $Y \rightarrow X$. ANLTSM and TS-LiNGAM only allow linear instantaneous effects. Theorem 1 shows that the causal summary time graph may still be identifiable when the instantaneous effects are linear and the variables are jointly Gaussian. TS-LiNGAM does not work in these situations. (2) *Confounders*: G-causality might fail when there is a confounder between X_t and Y_{t+1} , for example: The path between X_t and Y_{t+1} cannot be blocked by conditioning on any of the observed variables; G-causality infers $X \rightarrow Y$. ANLTSM does not allow for nonlinear confounders or confounders with time structure and TS-LiNGAM may fail, too (Exp. 1). (3) *Bad model assumptions*: The methods

share a similar problem: Performing general conditional independence tests is desirable, but not feasible, partially because the conditioning sets are too large [e.g. Bergsma, 2004]. Thus, the test is performed under a simple model, for example a linear one. If the model assumption is violated, one may draw wrong conclusions without noticing (e.g. Exp. 3). For TiMINo, that we define below, Lemma 1 shows that after fitting and checking the model by testing for independent residuals, the difficult conditional independences have been checked implicitly.

Thus, a *model check* is a simple but effective improvement. Although G-causality for two time series can easily be augmented with a cross-correlation test, we do not see a straight-forward extension to the multivariate G-causality. Furthermore, testing for cross-correlation does not always suffice (see Section 5.1).

4 Functional models for time series: TiMINo

We define TiMINo, a model class including the models described above and prove its identifiability.

Definition 1 Consider a time series $\mathbf{X}_t = (X_t^i)_{i \in V}$, such that the finite dimensional distributions are absolutely continuous with respect to a product measure (i.e. there is a pdf or a pmf). We say the time series satisfies a TiMINo if there is a $p > 0$ and if $\forall i \in V$ there are sets $\mathbf{PA}_0^i \subseteq X^{V \setminus \{i\}}$, $\mathbf{PA}_k^i \subseteq X^V$, s.t. $\forall t$

$$X_t^i = f_i((\mathbf{PA}_p^i)_{t-p}, \dots, (\mathbf{PA}_1^i)_{t-1}, (\mathbf{PA}_0^i)_t, N_t^i), \quad (1)$$

with N_t^i (jointly) independent and for each i , N_t^i identically distributed in t . The corresponding full time graph is obtained by drawing arrows from any node that appears in the right-hand side of (1) to X_t^i . We require the full time graph to be acyclic.

Below we assume that equations (1) follow an identifiable functional model class (IFMOC), Peters et al. [2011b] give a precise definition. Basically, it means that (I) *causal minimality* holds, a weak form of faithfulness that assumes a statistical dependence between cause and effect given all other parents [Spirtes et al., 2000]. And (II), all f_i come from a function class (e.g. additive noise) that is small enough to make the bivariate case identifiable (Section 2.2) if we exclude certain function-input-noise combinations like linear-Gaussian-Gaussian. The proof of the following theoretical result can be found in the appendix.

Theorem 1 Suppose that \mathbf{X}_t can be represented as a TiMINo with $\mathbf{PA}(X_t^i) = \bigcup_{k=0}^p (\mathbf{PA}_k^i)_{t-k}$ being the direct causes of X_t^i and that one of the following holds:

- (i) Equations (1) come from an IFMOC.
- (ii) Each component of the time series exhibits a time structure (i.e. $\mathbf{PA}(X_t^i)$ contains at least one X_{t-k}^i), the joint distribution is faithful with respect to the full time graph, and the summary time graph is acyclic.

Then the full time graph can be recovered from the joint distribution. In particular, the true causal summary time graph is identifiable. (Note that neither of the two conditions implies the other.)

Regarding (i): Many choices of a function class are possible [Peters et al., 2011b]. In practice, however, one still needs to fit those functions f_i from the data, which means for additive noise that estimating $\mathbf{E}[X_t^i | \mathbf{X}_{t-p}, \dots, \mathbf{X}_{t-1}]$ should be feasible. Different results show that stationarity and/or α mixing, or geometric ergodicity are required [e.g. Chu and Glymour, 2008]. In this work we consider VAR fitting: $f_i(p_1, \dots, p_r, n) = a_{i,1} \cdot p_1 + \dots + a_{i,r} \cdot p_r + n$, gam regression: $f_i(p_1, \dots, p_r, n) = f_{i,1}(p_1) + \dots + f_{i,r}(p_r) + n$ [e.g. Bell et al., 1996], and GP regression: $f_i(p_1, \dots, p_r, n) = f_i(p_1, \dots, p_r) + n$. Note that linear functions lead to the model of Hyvärinen et al. [2008] as a special case.

Regarding (ii): This condition nicely shows how the time structure does not only make the causal inference problem harder (the iid assumption is dropped), but also easier. In the iid case, for example, the true graph is not identifiable if all components are jointly Gaussian and the relationships are linear; with time structure it is. (TS-LiNGAM would fail, though.)

5 A practical method: TiMINo causality

The algorithm for TiMINo causality is based on the theoretical finding in Theorem 1. It takes the time series data as input and outputs either a DAG that estimates the summary time graph or remains undecided. In principle, it tries to fit a TiMINo model to the data and outputs the corresponding graph. If no model with independent residuals is found, it outputs “I do not know”. For a time series with many components, this gets intractable. In Section 6, we concentrate on time series without feedback loops, where we can exploit a more efficient method:

5.1 Full causal discovery

For additive noise models (ANMs) without time structure, Mooij et al. [2009] propose a procedure that recovers the structure without enumerating all possible DAGs. This procedure can be modified to be of use for time series (Algorithm 1). As reported by Mooij et al. [2009], the time complexity is $\mathcal{O}(d^2)$, where d is the number of time series, regarding fitting models and independence testing as atomic operations. To get the full time complexity, $\mathcal{O}(d^2)$ has to be multiplied by the sum of the complexity of the regression method and the independence test, both chosen by the user.

Algorithm 1 TiMINo causality

```

1: Input: Samples from a  $d$ -dimensional time series of length  $T$ :  $(\mathbf{X}_1, \dots, \mathbf{X}_T)$ , maximal order  $p$ 
2:  $S := (1, \dots, d)$ 
3: repeat
4:   for  $k$  in  $S$  do
5:     Fit TiMINo for  $X_t^k$  using  $X_{t-p}^k, \dots, X_{t-1}^k, X_{t-p}^i, \dots, X_{t-1}^i, X_t^i$  for  $i \in S \setminus \{k\}$ 
6:     Test if residuals are indep. of  $X^i, i \in S$ .
7:   end for
8:   Choose  $k^*$  to be the  $k$  with the weakest dependence. (If there is no  $k$  with independence, break and output: “I do not know - bad model fit”).
9:    $S := S \setminus \{k^*\}$ 
10:   $\text{pa}(k^*) := S$ 
11: until  $\text{length}(S)=1$ 
12: For all  $k$  remove all unnecessary parents.
13: Output:  $(\text{pa}(1), \dots, \text{pa}(d))$ 

```

Depending on the assumed model class, TiMINo causality has to be provided with a fitting method. Here, we chose `ar`, `gam` and `gptk` in R (<http://www.r-project.org/>) for linear models, generalized additive models, and GP regression. We call the methods TiMINo-linear, TiMINo-gam and TiMINo-GP, respectively. For the first two AIC determines the order of the process. All fitting methods are used in a “standard way”. For `gam` we used the built-in nonparametric smoothing splines. For the GP we used zero mean, squared exponential covariance function and Gaussian Likelihood. The hyper-parameters are automatically chosen by marginal likelihood optimization.

To test for independence between a residual time series N_t^k and another time series $X_t^i, i \in S$, we shift the latter time series up to the maximal order $\pm p$ (but at least up to ± 4); for each of those combinations we perform HSIC [Gretton et al., 2008], an independence test for iid data. One could also use a test based on cross-correlation that can be derived from Thm 11.2.3. in [Brockwell and Davis, 1991]. This is related to what is done in transfer function modeling [e.g. §13.1 in Brockwell and Davis, 1991], which is restricted to two time series and linear functions. But testing for cross-correlation is often not enough: if no time structure is present (iid data), it is obvious that correlation tests are most often insufficient. Also, experiments 1 and 5 describe situations, in which cross-correlations fail. To reduce the running time, however, one can use cross-correlation to determine the graph structure and use HSIC as a final model check. For HSIC we used a Gaussian kernel; as in [Gretton et al., 2008], the bandwidth is chosen to be the median distance of the input data. This is a heuristic but well-established choice.

Note that any other fitting method and independence test can be used as well. Although they work well in practice, we do not claim that our choices are optimal.

5.2 Partial causal discovery

Let \mathbf{X}_t “almost” satisfy a TiMINo model, that is some time series are unobserved or some functional relationships are not included in the model. We expect that the full discovery method remains undecided. One can modify the method such that it tries to discover parts of the causal graph: Whenever no k with independent residuals is found in line 8 of Algorithm 1 one subtracts a subset S_0 from the current version of S (first subtract one element, then any combination of two etc.) and repeat. If the method is able to fit a TiMINo model using only the remaining set $S \setminus S_0$, output this solution and S_0 , which has been excluded. Since there are $2^{\#S}$ subsets, this is only feasible for small S (see Exp. 6). This method may also be useful for the iid case; its theoretical properties remain to be investigated.

5.3 Weaknesses

(i) In principle, it may happen that the model assumption are violated, but one can nevertheless fit a model in the wrong direction (that is why we wrote “remaining *mostly* undecided”). This requires an “unnatural” fine tuning of the functions and Janzing and Steudel [2010] argue that in the case of causal sufficiency it cannot occur if one believes in the “independence” of cause and causal mechanism. Also, (i) is relevant only when there are time series without time structure or the data are non-faithful (see Theorem 1). We do not provide a precise analysis of the case with confounders, but analyze this situation empirically in Experiment 1. (ii) The null hypothesis of the independence test represents independence, although the scientific discovery of a causal relationship should rather be the alternative hypothesis. This fact may lead to wrong causal conclusions (instead of “I do not know”) on small data sets since we cannot reject independence for the wrong direction. This effect is strengthened by the Bonferroni correction of the HSIC based independence test. This may require modifications, when the number of time series is very high. It is thus useful to develop heuristics for “minimal” sample sizes. (iii) For large sample sizes, even smallest differences between the true data generating process and the model may lead to rejected independence tests [discussed by Peters et al., 2011a].

6 Experiments

Code is available in the suppl. mat. and will be online.

6.1 Artificial Data

We always included instantaneous effects, fitted models up to order $p = 2$ or $p = 6$ and set $\alpha = 0.05$.

Experiment 1: Confounder with time lag. We simulate 100 data sets (length 1000) from $Z_t = a \cdot Z_{t-1} + N_{Z,t}$, $X_t = 0.6 \cdot X_{t-1} + 0.5 \cdot Z_{t-1} + N_{X,t}$, $Y_t = 0.6 \cdot Y_{t-1} + 0.5 \cdot Z_{t-2} + N_{Y,t}$, with a between 0 and 0.95 and $N_{\cdot,t} \sim 0.4 \cdot \mathcal{N}(0, 1)^3$. Here, Z is a hidden common cause for X and Y . For all a , X_t contains information about Z_{t-1} and Y_{t+1} (see Figure 1); G-causality and TS-LiNGAM wrongly infer $X \rightarrow Y$. For large a , Y_t contains additional information about X_{t+1} , which leads to the wrong arrow $Y \rightarrow X$. TiMINo causality does not decide for any a . The nonlinear methods perform very similar (not shown). Note that for $a = 0$, a cross-correlation test is not enough to reject $X \rightarrow Y$. Further, all methods fail for $a = 0$ and Gaussian noise.

Experiment 2: Linear, Gaussian with instantaneous effects. We sample 100 data sets (length 2000) from $X_t = A_1 \cdot X_{t-1} + N_{X,t}$, $W_t = A_2 \cdot W_{t-1} + A_3 \cdot X_t + N_{W,t}$, $Y_t = A_4 \cdot Y_{t-1} + A_5 \cdot W_{t-1} + N_{Y,t}$, $Z_t = A_6 \cdot Z_{t-1} + A_7 \cdot W_t + A_8 \cdot Y_{t-1} + N_{Z,t}$ and $N_{\cdot,t} \sim 0.4 \cdot \mathcal{N}(0, 1)$ and A_i iid from $\mathcal{U}([-0.8, -0.2] \cup [0.2, 0.8])$. We regard the graph containing $X \rightarrow W \rightarrow Y \rightarrow Z$ and $W \rightarrow Z$ as correct. TS-LiNGAM and G-causality are not able to recover the true structure (see Table 1).

Experiment 3: Nonlinear, non-Gaussian without instantaneous effects. We simulate 100 data sets (length 500) from $X_t = 0.8X_{t-1} + 0.3N_{X,t}$, $Y_t = 0.4Y_{t-1} + (X_{t-1} - 1)^2 + 0.3N_{Y,t}$, $Z_t = 0.4Z_{t-1} +$

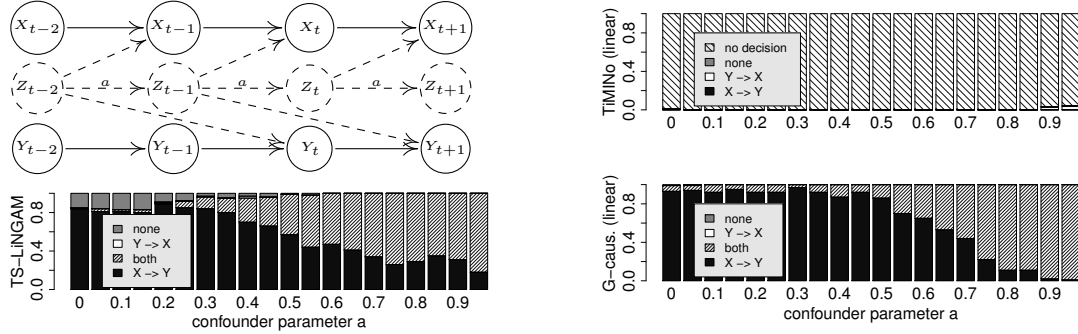


Figure 1: Exp.1: Part of the causal full time graph with hidden common cause Z (top left). TiMINo causality does not decide (top right), whereas G-causality and TS-LiNGAM wrongly infer causal connections between X and Y (bottom).

Table 1: Exp.2: Gaussian data and linear instantaneous effects: only TiMINo mostly discovers the correct DAG.

DAG	lin. Granger	TiMINo-lin	TS-LiNGAM
correct	13%	83%	19%
wrong	87%	7%	81%
no dec.	0%	10%	0%

$0.5 \cos(Y_{t-1}) + \sin(Y_{t-1}) + 0.3N_{Z,t}$, with $N_{.,t} \sim \mathcal{U}([-0.5, 0.5])$ (similar results for other noise distributions, e.g. exponential). Thus, $X \rightarrow Y \rightarrow Z$ is the ground truth. Nonlinear G-causality fails since the implementation is only pairwise and it thus always infers an effect from X to Z . Linear G-causality cannot remove the nonlinear effect from X_{t-2} to Z_t by using Y_{t-1} and gives many wrong answers. Also TiMINo-linear assumes a wrong model, but does not make any decision. TiMINo-gam and TiMINo-GP work well on this data set (Table 2). This specific choice of parameters show that a significant difference in performance is possible. For other parameters (e.g. less impact of the nonlinearity), G-causality and TS-LiNGAM still assume a wrong model but make fewer mistakes.

Experiment 4: Non-additive interaction. We simulate 100 data sets with different lengths from $X_t = 0.2 \cdot X_{t-1} + 0.9N_{X,t}$, $Y_t = -0.5 + \exp(-(X_{t-1} + X_{t-2})^2) + 0.1N_{Y,t}$, with $N_{.,t} \sim \mathcal{N}(0, 1)$. Figure 2 shows that TiMINo-linear and TiMINo-gam remain mainly undecided, whereas TiMINo-GP performs well. For small sample sizes, one observes two effects: GP regression does not obtain accurate estimates for the residuals, these estimates are not independent and thus TiMINo-GP remains more often undecided. Also, TiMINo-gam makes more correct answers than one would expect due to more type II errors. Linear G-causality and TS-LiNGAM give more than 90% incorrect answers, but non-linear G-causality is most often correct (not shown). Bad model assumptions do not *always* lead to incorrect causal conclusions.

Experiment 5: Non-linear Dependence of Residuals. In Experiment 1, TiMINo equipped with a cross-correlation inferred a causal edge, although there were none. The opposite is also possible: $X_t = -0.5 \cdot X_{t-1} + N_{X,t}$, $Y_t = -0.5 \cdot Y_{t-1} + X_{t-1}^2 + N_{Y,t}$ and $N_{.,t} \sim 0.4 \cdot \mathcal{N}(0, 1)$ (length 1000). TiMINo-gam with cross-correlation infers no causal link between X and Y , whereas TiMINo-gam with HSIC correctly identifies $X \rightarrow Y$.

Experiment 6: Partial Causal Discovery. We sample 100 data sets (length 600) from $X_t = 0.5 \cdot X_{t-1} + N_{X,t}$, $B_t = 0.5 \cdot B_{t-1} + N_{B,t}$, $A_t = 0.5 \cdot A_{t-1} + 0.5 \cdot B_{t-1} + N_{A,t}$, $Y_t = 0.5 \cdot Y_{t-1} - 0.9 \cdot X_{t-1} + 0.8 \cdot B_{t-1} + N_{Y,t}$, $W_t = 0.5 \cdot W_{t-1} + 0.8 \cdot X_{t-1} + N_{W,t}$ and $N_{.,t} \sim 0.4 \cdot \mathcal{U}([-0.5, 0.5])$. Let X_t be latent. The standard method finds A_t as a “sink time series” and halts in iteration two (line 8 in Algorithm 1). Instead of outputting “I do not know”, the partial discovery method described in Section 5.2 is able to correctly infer this DAG (see Figure 3) in 82% of the cases (18% wrong answers). G-causality and

Table 2: Exp.3: Since the data are nonlinear, linear G-causality and TS-LiNGAM give wrong answers, TiMINo-lin does not decide. Nonlinear G-causality fails because it analyzes the causal structure between pairs of time series.

DAG	Granger (lin)	Granger (nonlin)	TiMINo (lin)	TiMINo (gam)	TiMINo (GP)	TS-LiNGAM
correct	69%	0%	0%	95%	94%	12%
wrong	31%	100%	0%	1%	1%	88%
no dec.	0%	0%	100%	4%	5%	0%

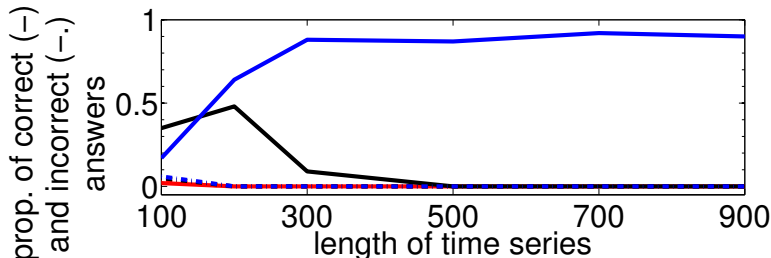


Figure 2: Exp.4: TiMINo-GP (blue) works reliably for long time series. TiMINo-linear (red) and TiMINo-gam (black) mostly remain undecided.

TS-LiNGAM give only wrong answers.

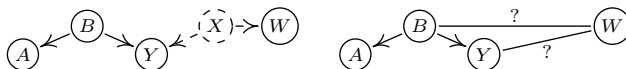


Figure 3: Exp.6: The true causal summary time graph (left) cannot be recovered because X_t is unobserved. TiMINo gives a partial result (right).

6.2 Real Data

We fitted up to order 6 and included instantaneous effects. For TiMINo, “correct” means that TiMINo-gam makes the correct decision and TiMINo-linear is correct or undecided. TiMINo-GP always remains undecided because there are too few data points to fit such a general model. Again, α is set to 0.05.

Experiment 7: Gas Furnace. [Box et al., 2008, length 296], X_t describes the input gas rate and Y_t the output CO₂. We regard $X \rightarrow Y$ as being true. TS-LiNGAM, G-causality, TiMINo-lin and TiMINo-gam correctly infer $X \rightarrow Y$. Disregarding time information leads to a wrong causal conclusion: The method described by Hoyer et al. [2009] leads to a p -value of 4.8% in the correct and 9.1% in the false direction.

Experiment 8: Old Faithful. [Azzalini and Bowman, 1990, length 194] X_t contains the duration of an eruption and Y_t the time interval to the next eruption of the Old Faithful geyser. We regard $X \rightarrow Y$ as the ground truth. Although the time intervals $[t, t + 1]$ do not have the same length for all t , we model the data as two time series. TS-LiNGAM and TiMINo give correct answers, whereas linear G-causality infers $X \rightarrow Y$, and nonlinear G-causality infers $Y \rightarrow X$.

Experiment 9: Temperature. (available at <https://webdav.tuebingen.mpg.de/cause-effect/>, length 16382) X_t are indoor and Y_t outdoor measurements (recorded every 5 minutes), we expect that there is a causal link $Y \rightarrow X$. TS-LiNGAM wrongly infers $X \rightarrow Y$ and both G-causality methods infer a bidirected arrow. TiMINo remains undecided. Maybe, the data are causal insufficient: time may con-

found outdoor temperature and the usage of heating, the latter is a direct cause for indoor temperature. Also, Y may cause heating. Such a model does not allow for a TiMINo from Y to X .

Experiment 10: Abalone (no time structure). The abalone data set [Asuncion and Newman, 2007] contains (among others that lead to similar results) age X_t and diameter Y_t of a certain shell fish. If we model 1000 randomly chosen samples as time series, G-causality (both linear and nonlinear) infers no causal relation as expected. TS-LiNGAM wrongly infers $Y \rightarrow X$, which is probably due to the nonlinear relationship. TiMINo gives the correct result.

Experiment 11: Diary (confounder). We consider 10 years of weekly prices for butter X_t and cheddar cheese Y_t [Gould, 2007, length 522]. They are strongly correlated, but we expect this correlation to be due to the (hidden) milk price M_t : $X \leftarrow M \rightarrow Y$. TiMINo does not decide, whereas TS-LiNGAM and G-causality wrongly infer $X \rightarrow Y$. This may be due to different time lags of the confounder (cheese has longer storing and maturing times than butter).

The phase slope index [Nolte et al., 2008] performed well only in Exp. 6, in all other experiments it either gave wrong results or did not decide. Due to space constraints we omit details about this method.

7 Conclusions and Future Work

This paper shows how causal inference benefits from the framework of functional models. TiMINo causality can be seen as an extension of methods from the iid case, but the benefits compared to other time series methods are substantial and important: It comes with an identifiability that is more general than existing results and lead to a practical algorithm that allows for the ability to make no decision instead of a wrong one. TiMINo is applicable to multivariate, linear, nonlinear and instantaneous interactions and can also discover partial structures. On the data sets considered it outperforms existing methods.

We think the following investigations would be worthwhile: (1) Applying more complex models (like heteroscedastic models) and preprocessing the data (removing trends, periodicities, etc.) may decrease the number of cases where TiMINo causality is undecided. (2) Checking for autocorrelations in the residuals is another possible model check and not included yet. (3) In the case of non-instantaneous feedback loops, one should find a method to fit the model structure that is faster than brute-force search. (4) Although we report promising results, an extensive evaluation of this method on even more real data sets is necessary. This lies beyond the scope of the present conference paper.

8 Appendix

Lemma 1 *If $\mathbf{X}_t = (X_t^i)_{i \in V}$ satisfy a TiMINo model, each variable X_t^i is conditionally independent of each of its non-descendants given its parents.*

Proof . With $\mathcal{S} := \mathbf{PA}(X_t^i) = \bigcup_{k=0}^p (\mathbf{PA}_k^i)_{t-k}$ and equation (1) we get $X_t^i |_{\mathcal{S}=s} = f_i(s, N_t^i)$ for an s with $p(s) > 0$. Any non-descendant of X_t^i can be written as a function of all noise variables from its ancestors and is therefore independent of X_t^i given $\mathcal{S} = s$. For this proof it is crucial that we consider time series for $t \in \mathbb{N}$. We believe that a similar statement holds for $t \in \mathbb{Z}$, which only introduces technical difficulties. \square

Proof of Theorem 1 Suppose that \mathbf{X}_t allows two different representations of TiMINo that lead to two different full time graphs \mathcal{G} and \mathcal{G}' . (i) First we assume that \mathcal{G} and \mathcal{G}' do not differ in the instantaneous effects: $\mathbf{PA}_0^i(\text{in } \mathcal{G}) = \mathbf{PA}_0^i(\text{in } \mathcal{G}') \forall i$. Without loss of generality, there is some $k > 0$ and an edge $X_{t-k}^1 \rightarrow X_t^2$, say, that is in \mathcal{G} but not in \mathcal{G}' . From \mathcal{G}' and Lemma 1 we have that $X_{t-k}^1 \perp\!\!\!\perp X_t^2 | \mathcal{S}$, where $\mathcal{S} = (\{X_{t-l}^i, 1 \leq l \leq p, i \in V\} \cup \mathbf{ND}_t) \setminus \{X_{t-k}^1, X_t^2\}$, and \mathbf{ND}_t are all X_t^i that are non-descendants (wrt instantaneous effects) of X_t^2 . Applied to \mathcal{G} , causal minimality leads to a contradiction: $X_{t-k}^1 \not\perp\!\!\!\perp X_t^2 | \mathcal{S}$. Now we suppose \mathcal{G} and \mathcal{G}' differ in the instantaneous effects. This time we choose $\mathcal{S} = \{X_{t-l}^i, 1 \leq l \leq p, i \in V\}$. Then for each s and i we have: $X_t^i |_{\mathcal{S}=s} = f_i(s, (\tilde{\mathbf{PA}}_0^i)_t)$, where $\tilde{\mathbf{PA}}_0^i$ are all instantaneous parents of X_t^i conditioned on $\mathcal{S} = s$. All $X_t^i |_{\mathcal{S}=s}$ with the instantaneous effects describe two different structures of an IFMOC. This contradicts the identifiability results by Peters et al. [2011b]. (ii) Because

of Lemma 1 and faithfulness \mathcal{G} and \mathcal{G}' only differs in the instantaneous effects. But each instantaneous arrow $X_t^i \rightarrow X_t^j$ forms a v -structure together with $X_{t-k}^j \rightarrow X_t^j$; the latter exists because of the time structure and X_{t-k}^j cannot be connected with X_t^i since this introduces a cycle in the summary time graph. \square

References

- N. Ancona, D. Marinazzo, and S. Stramaglia. Radial basis function approach to nonlinear granger causality of time series. *Phys. Rev. E*, 70(5):056221, 2004.
- A. Asuncion and D. J. Newman. UCI repository. <http://archive.ics.uci.edu/ml/>, 2007.
- A. Azzalini and A. W. Bowman. A look at some data on the Old Faithful Geyser. *Applied Statistics*, 39(3):357–365, 1990.
- D. Bell, J. Kay, and J. Malley. A non-parametric approach to non-linear causality testing. *Economics Letters*, 51(1):7–18, 1996.
- W. P. Bergsma. *Testing conditional independence for continuous random variables*, 2004. EURANDOM-report 2004-049.
- G. E. P. Box, G. M. Jenkins, and G. C. Reinsel. *Time series analysis: forecasting and control*. Wiley series in probability and statistics. John Wiley, 2008.
- P. J. Brockwell and R. A. Davis. *Time Series: Theory and Methods*. Springer, 2nd edition, 1991.
- Y. Chen, G. Rangarajan, J. Feng, and M. Ding. Analyzing multiple nonlinear time series with extended granger causality. *Physics Letters A*, 324, 2004.
- T. Chu and C. Glymour. Search for additive nonlinear time series causal models. *Journal of Machine Learning Research*, 9:967–991, 2008.
- M. Eichler. Graphical modelling of multivariate time series. *Probability Theory and Related Fields*, pages 1–36, 2011.
- J. P. Florens and M. Mouchart. A note on noncausality. *Econometrica*, 50(3):583–591, 1982.
- Brian W. Gould. Diary data sets. <http://future.aae.wisc.edu/tab/prices.html>, 2007.
- C. W. J. Granger. Investigating causal relations by econometric models and cross-spectral methods. *Econometrica*, 37(3):424–38, July 1969.
- A. Gretton, K. Fukumizu, C. H. Teo, L. Song, B. Schölkopf, and A. Smola. A kernel statistical test of independence. In *NIPS 20*, Canada, 2008.
- T. J. Hastie and R. J. Tibshirani. *Generalized additive models*. London: Chapman & Hall, 1990.
- P. Hoyer, D. Janzing, J. Mooij, J. Peters, and B. Schölkopf. Nonlinear causal discovery with additive noise models. In *NIPS 21*, Canada, 2009.
- A. Hyvärinen, S. Shimizu, and P. O. Hoyer. Causal modelling combining instantaneous and lagged effects: an identifiable model based on non-gaussianity. In *ICML 25*, 2008.
- D. Janzing and B. Steudel. Justifying additive-noise-model based causal discovery via algorithmic information theory. *Open Systems and Information Dynamics*, 17:189–212, 2010.
- J. Mooij, D. Janzing, J. Peters, and B. Schölkopf. Regression by dependence minimization and its application to causal inference. In *ICML 26*, 2009.

- G. Nolte, A. Ziehe, V. Nikulin, A. Schlögl, N. Krämer, T. Brismar, and K.-R. Müller. Robustly Estimating the Flow Direction of Information in Complex Physical Systems. *Phys. Rev. Letters*, 100, 2008.
- J. Pearl. *Causality: Models, reasoning, and inference*. Cambridge University Press, 2nd edition, 2009.
- J. Peters, D. Janzing, and B. Schölkopf. Causal inference on discrete data using additive noise models. *IEEE Trans. Pattern Analysis Machine Intelligence*, 33(12):2436–2450, 2011a.
- J. Peters, J. M. Mooij, D. Janzing, and B. Schölkopf. Identifiability of causal graphs using functional models. In *UAI 27*, 2011b.
- S. Shimizu, P. O. Hoyer, A. Hyvärinen, and A. J. Kerminen. A linear non-Gaussian acyclic model for causal discovery. *Journal of Machine Learning Research*, 7:2003–2030, 2006.
- P. Spirtes, C. Glymour, and R. Scheines. *Causation, Prediction, and Search*. MIT Press, 2nd edition, 2000.
- K. Zhang and A. Hyvarinen. On the identifiability of the post-nonlinear causal model. In *UAI 25*, 2009.