

A constraint optimization approach to causal discovery from subsampled time series data [☆]



Antti Hyttinen ^{a,*}, Sergey Plis ^b, Matti Järvisalo ^a, Frederick Eberhardt ^c,
David Danks ^d

^a *HIT, Department of Computer Science, University of Helsinki, Finland*

^b *Mind Research Network and University of New Mexico, United States*

^c *Humanities and Social Sciences, California Institute of Technology, United States*

^d *Department of Philosophy, Carnegie Mellon University, United States*

ARTICLE INFO

Article history:

Received 1 December 2016

Received in revised form 30 June 2017

Accepted 10 July 2017

Available online 29 July 2017

Keywords:

Causality

Causal discovery

Graphical models

Time series

Constraint satisfaction

Constraint optimization

ABSTRACT

We consider causal structure estimation from time series data in which measurements are obtained at a coarser timescale than the causal timescale of the underlying system. Previous work has shown that such subsampling can lead to significant errors about the system's causal structure if not properly taken into account. In this paper, we first consider the search for system timescale causal structures that correspond to a given measurement timescale structure. We provide a constraint satisfaction procedure whose computational performance is several orders of magnitude better than previous approaches. We then consider finite-sample data as input, and propose the first constraint optimization approach for recovering system timescale causal structure. This algorithm optimally recovers from possible conflicts due to statistical errors. We then apply the method to real-world data, investigate the robustness and scalability of our method, consider further approaches to reduce underdetermination in the output, and perform an extensive comparison between different solvers on this inference problem. Overall, these advances build towards a full understanding of non-parametric estimation of system timescale causal structures from subsampled time series data.

© 2017 Elsevier Inc. All rights reserved.

1. Introduction

Time-series data has long constituted the basis for causal modeling in many fields of science [12,15,22]. These data often provide very precise measurements at regular time points, but the underlying causal interactions that give rise to those measurements can occur at a much faster timescale than the measurement frequency. As just one example: fMRI experiments measure neural activity (given various assumptions) roughly once per two seconds, but the underlying neural connections clearly operate much more quickly. Time order information can simplify causal analysis since it can provide

[☆] This paper is part of the Virtual special issue on the Eighth International Conference on Probabilistic Graphical Models, Edited by Giorgio Corani, Alessandro Antonucci, Cassio De Campos.

* Corresponding author.

E-mail addresses: antti.hyttinen@helsinki.fi (A. Hyttinen), s.m.plis@gmail.com (S. Plis), matti.jarvisalo@helsinki.fi (M. Järvisalo), fde@caltech.edu (F. Eberhardt), ddanks@cmu.edu (D. Danks).

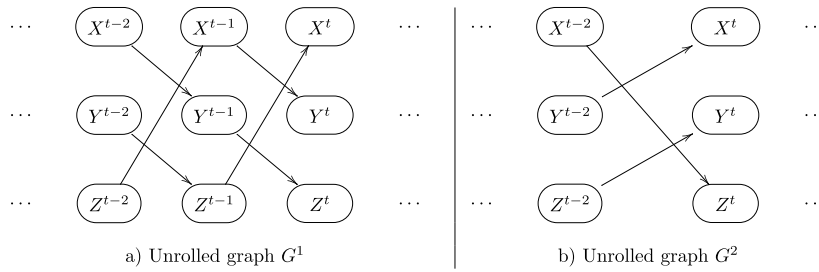


Fig. 1. (a) The structure of the causal system-scale time series. (b) The structure of the corresponding measurement scale time series if only every second sample is observed i.e. nodes at time slice $t - 1$ are marginalized. If subsampling is ignored and (b) is thought to depict the true causal structure, all direct causal relationships among $\{X, Y, Z\}$ are misspecified.

directionality, but time series data that undersamples the generating process can be especially misleading about the true direct causal connections [7,19].

For example, Fig. 1a shows the causal structure of a process unrolled over discrete time steps, and Fig. 1b shows the corresponding structure of the same process, obtained by marginalizing every second time step. If we do not take into account the possibility of subsampling, then we would conclude that Fig. 1b gives the correct structure – and thus totally miss the presences of all true edges. This drastic structure misspecification may lead us to perform a possibly costly intervention on Z to control Y , when the influence of Z on Y is, in fact, completely mediated by X and so, intervening on X would be a more effective choice. Also, a (parametric) model with the structure in Fig. 1b gives inaccurate predictions when intervening on both X and Z : the value of Y would be predicted to depend on Z and not on X , when in reality Y depends on X and not on Z .

Standard methods for estimating causal structure from time series either focus exclusively on estimating a transition model at the measurement timescale (e.g., Granger causality [12,13]) or combine a model of measurement timescale transitions with so-called “instantaneous” or “contemporaneous” causal relations that aim to capture interactions that are faster than the measurement process (e.g., SVAR [22,15,18]), though only very specific types of interactions can be captured with these latter models. In contrast, we follow Plis et al. [30,31] and Gong et al. [11], and explore the possibility of identifying (features of) the causal process at the true timescale from data that subsample this process.

Plis et al. [30,31] developed algorithms that can learn the set of causal timescale structures that could yield a given measurement timescale graph, either at a known or unknown undersampling rate. While these algorithms show that the inference problem is solvable, they face a number of computational challenges that limit their use. They do, however, show the importance of constraints for this problem, and so suggest that a constraint satisfaction approach might be more effective and efficient. Gong et al. [11] consider finding a linear SVAR from subsampled data. They show that if the error variables are non-Gaussian, the true causal effects matrix can be discovered even from subsampled data. However, their method is highly restricted in terms of numbers of variables and parametric form.

In this paper, we provide an exact discovery algorithm based on using a general-purpose Boolean constraint solver [4, 10], and demonstrate that it is orders of magnitudes faster than the current state-of-the-art method by Plis et al. [31]. At the same time, our approach is much simpler and, as we show, it allows inference in more general settings. We then develop the approach to integrate possibly conflicting constraints obtained from the data. In addition to an application of the method to the real-world data, we investigate the robustness and scalability of our method, consider further approaches to reduce underdetermination in the output, and perform an extensive comparison between different solvers on this inference problem. Moreover, unlike the method by Gong et al. [11], our approach does not depend on a particular parameterization of the underlying model and scales to a more reasonable number of variables.

The code implementing the approach presented in this article, including the answer set programming and Boolean satisfiability encodings, is available at

<http://www.cs.helsinki.fi/group/coreo/subsampled/>.

This article considerably extends a preliminary version presented at the International Conference on Probabilistic Graphical Models 2016 (PGM 2016) [17]. Most noticeably, Sections 6–9 of this article provide entirely new contents, including a real-world case study (Section 6), an evaluation of the impact of the choice of constraint satisfaction and optimization solvers on the efficiency of the approach (Section 7), and a discussion on learning from mixed frequency data (Section 8). Furthermore, new simulations on accuracy and robustness (Section 5, Figures 7–9) are now included.

2. Representation

We assume that the system of interest relates a set of variables $\mathbf{V}^t = \{X^t, Y^t, Z^t, \dots\}$ defined at discrete time points $t \in \mathbb{Z}$ with continuous ($\in \mathbb{R}^n$) or discrete ($\in \mathbb{Z}^n$) values [9]. We distinguish the representation of the true causal process at the *system or causal timescale* from the time series data that are obtained at the *measurement timescale*. Following Plis et al.

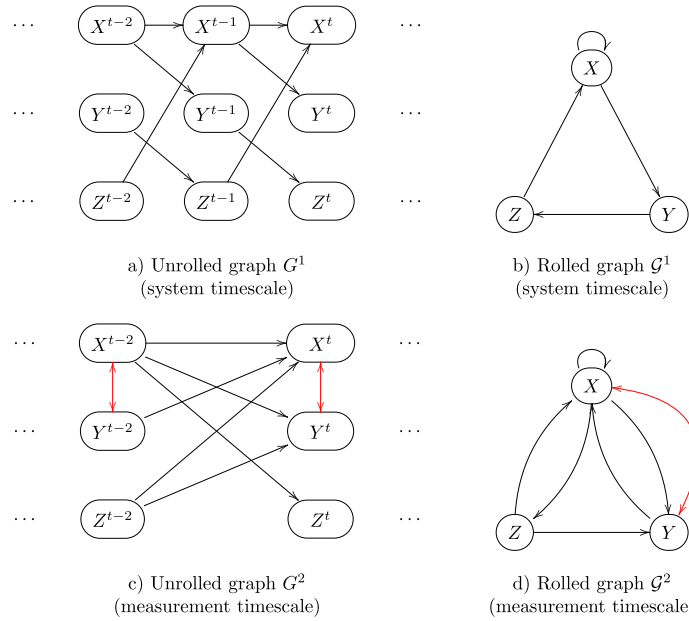


Fig. 2. Graph (a) shows the unrolled system timescale structure, where edges repeat through time steps. Graph (b) shows the rolled representation of the same structural information. Graph (c) shows the measurement timescale structure for subsampling rate $u = 2$, i.e. nodes at time slice $t - 1$ in graph (a) are marginalized. Graph (d) depicts the rolled representation of the same structural information as in graph (c).

[31], we assume that the true between-variable causal interactions at the system timescale constitute a first-order Markov process; that is, that the independence $\mathbf{V}^t \perp\!\!\!\perp \mathbf{V}^{t-k} | \mathbf{V}^{t-1}$ holds for all $k > 1$. The parametric models for these causal structures are structural vector autoregressive (SVAR) processes or dynamic (discrete/continuous variable) Bayes nets. Since the system timescale can be arbitrarily fast (and causal influences take time), we assume that there is no “contemporaneous” causation of the form $X^t \rightarrow Y^t$ [14]. We also assume that \mathbf{V}^{t-1} contains all common causes of variables in \mathbf{V}^t . These assumptions jointly express the widely used causal sufficiency assumption (see [35]) in the time series setting. In this non-parametric setting, we consider surgical interventions (on the observed variables in \mathbf{V}) that keep variables fixed at the selected values through the (causal timescale) time steps.

The system timescale causal structure can thus be represented by a causal graph G^1 (as in a dynamic Bayes net) with edges only of the form $X^{t-1} \rightarrow Y^t$, where $X = Y$ is permitted (see Fig. 2a for an example). Since the causal process is time-invariant, the edges repeat through t . In accordance with Plis et al. [31], for any G^1 we use a simpler, rolled graph representation, denoted by \mathcal{G}^1 , where for all $X, Y: X \rightarrow Y \in \mathcal{G}^1$ iff $X^{t-1} \rightarrow Y^t \in G^1$. That is, the rolled graph represents time only implicitly in the edges, rather than through variable duplication. Both the unrolled and rolled representations contain exactly the same structural information. Fig. 2b shows the rolled graph representation \mathcal{G}^1 of G^1 in Fig. 2a.

Time series data are obtained from the above process at the *measurement timescale*, defined by some (possibly unknown) integral sampling rate u . The measured time series sample \mathbf{V}^t is at times $t, t - u, t - 2u, \dots$; we are interested in the case of $u > 1$, i.e., the case of subsampled data. A different route to subsampling would use continuous-time models as the underlying system timescale structure. However, some series (e.g., transactions such as salary payments) are inherently discrete-time processes [11], and many continuous-time systems can be approximated arbitrarily closely as discrete-time processes. Thus, we focus here on discrete-time causal structures as a justifiable, yet simple, basis for our non-parametric inference procedure.

The (causal) structure of this subsampled time series can be obtained (leaving aside sampling variation) from G^1 by marginalizing the intermediate time steps. Fig. 2c shows the measurement timescale structure G^2 corresponding to subsampling rate $u = 2$ for the system timescale causal structure in Fig. 2a. Each directed edge in G^2 corresponds to a directed path of length 2 in G^1 . For arbitrary u, X, Y , the formal relationship between G^u and G^1 edges is

$$X^{t-u} \rightarrow Y^t \in G^u \Leftrightarrow X^{t-u} \rightsquigarrow Y^t \in G^1,$$

where \rightsquigarrow denotes a directed path.

G^u must also represent “direct” connections between variables in the same time step [37]. The bi-directed arrow $X^t \leftrightarrow Y^t$ in Fig. 2c is an example: X^{t-1} is an unobserved (in the data) common cause of X^t and Y^t in G^1 (Fig. 2a). Formally, the system timescale structure G^1 induces bi-directed edges in the measurement timescale G^u as follows:

$$X^t \leftrightarrow Y^t \in G^u \Leftrightarrow \exists Z, l < u : (X^t \leftarrow Z^{t-l} \rightsquigarrow Y^t) \in G^1, \quad \text{where } X \neq Y.$$

Just as \mathcal{G}^1 represents the rolled version of G^1 , \mathcal{G}^u represents the rolled version of G^u : $X \rightarrow Y \in \mathcal{G}^u$ iff $X^{t-u} \rightarrow Y^t \in G^u$ and $X \leftrightarrow Y \in \mathcal{G}^u$ iff $X^t \leftrightarrow Y^t \in G^u$.

The relationship between \mathcal{G}^1 and \mathcal{G}^u —that is, the impact of subsampling—can be concisely represented using only the rolled graphs:

$$X \rightarrow Y \in \mathcal{G}^u \Leftrightarrow X \overset{u}{\rightsquigarrow} Y \in \mathcal{G}^1, \tag{1}$$

$$X \leftrightarrow Y \in \mathcal{G}^u \Leftrightarrow \exists Z, l < u : (X \overset{l}{\rightsquigarrow} Z \overset{l}{\rightsquigarrow} Y) \in \mathcal{G}^1, \quad \text{where } X \neq Y. \tag{2}$$

Here $\overset{l}{\rightsquigarrow}$ denotes a path of length l . Using the rolled graph notation, the logical encodings in Section 3 are considerably simpler.

Subsampling can also be interpreted as a transitive operation applied to graphs. For example, \mathcal{G}^6 is the graph that results from subsampling \mathcal{G}^2 by a further factor of 3. More generally, $\mathcal{G}^{u \cdot k}$ can be obtained by subsampling \mathcal{G}^k by (another) u steps according to:

$$X \rightarrow Y \in \mathcal{G}^{u \cdot k} \Leftrightarrow X \overset{u}{\rightsquigarrow} Y \in \mathcal{G}^k,$$

$$X \leftrightarrow Y \in \mathcal{G}^{u \cdot k} \Leftrightarrow \exists Z, l < u : (X \overset{l}{\rightsquigarrow} Z \overset{l}{\rightsquigarrow} Y) \in \mathcal{G}^k \quad \vee$$

$$\exists Z, W, l < u : (X \overset{l}{\rightsquigarrow} Z \leftrightarrow W \overset{l}{\rightsquigarrow} Y) \in \mathcal{G}^k, \quad \text{where } X \neq Y.$$

Notice that in the latter equation, the bidirected edges in \mathcal{G}^k may induce additional bidirected edges in $\mathcal{G}^{u \cdot k}$. These equations yield Equations (1) and (2) when $k = 1$, since there are no bidirected edges in \mathcal{G}^1 .

In order to obtain a correspondence between the underlying causal structure and the distribution that gives rise to the observed data at measurement timescale, we assume for a given subsampling rate u that specific conditional independences correspond to the absence of specific causal connections:

$$X^{t-u} \perp\!\!\!\perp Y^t \mid \mathbf{V}^{t-u} \setminus X^{t-u} \Leftrightarrow X \rightarrow Y \notin \mathcal{G}^u \tag{3}$$

$$X^t \perp\!\!\!\perp Y^t \mid \mathbf{V}^{t-u} \Leftrightarrow X \leftrightarrow Y \notin \mathcal{G}^u \tag{4}$$

These assumptions are analogous to the combination of the Markov and faithfulness assumptions in the standard setting of causal discovery from cross-sectional data. However, here the assumptions are restricted to the particular (in)dependence relations we require to determine the causal structure, i.e., we allow, for example, for canceling pathways, which would otherwise constitute a violation of faithfulness, at subsampling rates that we do not consider.

Danks and Plis [6] demonstrated that, in the infinite sample limit, the causal structure \mathcal{G}^1 at the system timescale is in general underdetermined, even when the subsampling rate u is known and small. Consequently, even when ignoring estimation errors, the most we can learn is an equivalence class of causal structures at the system timescale. We define \mathcal{H} to be the estimated version of \mathcal{G}^u , a graph over \mathbf{V} obtained or estimated at the measurement timescale (with possibly unknown u). Due to underdetermination, multiple (\mathcal{G}^1, u) pairs can imply \mathcal{H} , and so search is particularly challenging when u is unknown. At the same time, if \mathcal{H} is estimated from data, it is possible, due to statistical errors, that no \mathcal{G}^u has the same structure as \mathcal{H} . With these observations, we are ready to define the computational problems focused on in this work.

Task 1. Given a measurement timescale structure \mathcal{H} (with possibly unknown u), infer the (equivalence class of) causal structures \mathcal{G}^1 consistent with \mathcal{H} (i.e. $\mathcal{G}^u = \mathcal{H}$ by Eqs. (1) and (2)) if such a \mathcal{G}^1 exists.

We also consider the corresponding problem when the subsampled time series is directly provided as input, rather than \mathcal{G}^u .

Task 2. Given a dataset of measurements of \mathbf{V} obtained at the measurement timescale (with possibly unknown u), infer the (equivalence class of) causal structures \mathcal{G}^1 (at the system timescale) that are (optimally) consistent with the data.

Section 3 provides a solution to Task 1. Section 4 provides a solution to Task 2, including an explanation on how \mathcal{H} can be estimated from sample data in Section 4.2. Later sections further consider generalizations of these two basic tasks.

3. Finding consistent system timescale structures

We first focus on Task 1. We discuss the computational complexity of the underlying decision problem, and present a practical Boolean constraint satisfaction approach that empirically scales up to significantly larger graphs than previous state-of-the-art algorithms.

3.1. On computational complexity

Consider the task of finding even a single \mathcal{G}^1 consistent with a given \mathcal{H} . A variant of the associated decision problem is related to the NP-complete problem of finding a matrix root.

Theorem 1. *Deciding whether there is a \mathcal{G}^1 that is consistent with the directed edges of a given \mathcal{H} is NP-complete for any fixed $u \geq 2$.*

Proof. Membership in NP follows from a guess and check: guess a candidate \mathcal{G}^1 , and deterministically check whether the length- u paths of \mathcal{G}^1 correspond to the edges of \mathcal{H} [31]. For NP-hardness, for any fixed $u \geq 2$, there is a straightforward reduction from the NP-complete problem of determining whether a Boolean B matrix¹ has a u th root [21]: for a given $n \times n$ Boolean matrix B , interpret B as the directed edge relation of \mathcal{H} , i.e., \mathcal{H} has the edge (i, j) iff $A^u(i, j) = 1$. It is then easy to see that there is a \mathcal{G}^1 that is consistent with the obtained \mathcal{H} iff $B = A^u$ for some binary matrix A (i.e., a u th root of B). \square

If u is unknown, then membership in NP can be established in the same way by guessing both a candidate \mathcal{G}^1 and a value for u . Theorem 1 ignores the possible bi-directed edges in \mathcal{H} (whose presence/absence is also harder to determine reliably from practical numbers of samples; see Section 5). Knowledge of the presences and absences of such edges in \mathcal{H} can restrict the set of candidate \mathcal{G}^1 s. For example, in the special case where \mathcal{H} is known to not contain any bi-directed edges, the possible \mathcal{G}^1 s have a fairly simple structure: in any \mathcal{G}^1 that is consistent with \mathcal{H} , every node has at most one successor.² Whether this knowledge can be used to prove a more fine-grained complexity result for special cases is an open question.

3.2. A SAT-based approach

Recently, the first exact search algorithm for finding the \mathcal{G}^1 s that are consistent with a given \mathcal{H} for a known u was presented by Plis et al. [31]; it represents the current state of the art. Their approach implements a specialized depth-first search procedure for the problem, with domain-specific polynomial time search-space pruning techniques. As an alternative, we present here a Boolean satisfiability based approach. First, we represent the problem exactly using a rule-based constraint satisfaction formalism. Then, for a given input \mathcal{H} , we employ an off-the-shelf Boolean constraint satisfaction solver for finding a \mathcal{G}^1 that is guaranteed to be consistent with \mathcal{H} (if such \mathcal{G}^1 exists). Our approach is not only simpler than the approach of Plis et al. [31], but as we will show, it also significantly improves the current state-of-the-art in runtime efficiency and scalability.

We present our approach using answer set programming (ASP) as the constraint satisfaction formalism³ [28,33,10]. It offers an expressive declarative modeling language, in terms of first-order logical rules, for various types of NP-hard search and optimization problems. To solve a problem via ASP, one first needs to develop an ASP program (in terms of ASP rules/constraints) that models the problem at hand; that is, the declarative rules implicitly represent the set of solutions to the problem in a precise fashion. Then one or multiple (optimal, in case of optimization problems) solutions to the original problem can be obtained by invoking an off-the-shelf ASP solver, such as the state-of-the-art Clingo system [10] used in this work. The search algorithms implemented in the Clingo system are extensions of state-of-the-art Boolean satisfiability and optimization techniques which can today outperform even specialized domain-specific algorithms, as we show here.

We proceed by describing a simple ASP encoding of the problem of finding a \mathcal{G}^1 that is consistent with a given \mathcal{H} . The input—the measurement timescale structure \mathcal{H} —is represented as follows. The input predicate `node/1` represents the nodes of \mathcal{H} (and all graphs), indexed by $1 \dots n$. The presence of a directed edge $X \rightarrow Y$ between nodes X and Y is represented using the predicate `edgeh/2` as `edgeh(X, Y)`. Similarly, the fact that an edge $X \rightarrow Y$ is not present is represented using the predicate `no_edgeh/2` as `no_edgeh(X, Y)`. The presence of a bidirected edge $X \leftrightarrow Y$ between nodes X and Y is represented using the predicate `confh/2` as `confh(X, Y)` ($X < Y$), and the fact that an edge $X \leftrightarrow Y$ is not present is represented using the predicate `no_confh/2` as `no_confh(X, Y)`.

If u is known, then it can be passed as input using `u(U)`; alternatively, it can be defined as a single value in a given range (here set to $1, \dots, 5$ as an example):

```
urange(1..5). % Define a range of u:s
1 { u(U) : urange(U) } 1. % u(U) is true for only one U in the range
```

Here the *cardinality constraint* `1 { u(U) : urange(U) } 1` states that the predicate `u` is true for exactly one value `U` chosen from those for which `urange(U)` is true.

¹ Multiplication of two values in $\{0, 1\}$ is defined as the logical-or, or equivalently, the maximum operator.

² To see this, assume X has two successors, Y and Z , s.t. $Y \neq Z$ in \mathcal{G}^1 . Then \mathcal{G}^u will contain a bi-directed edge $Y \leftrightarrow Z$ for all $u \geq 2$, which contradicts the assumption that \mathcal{H} has no bi-directed edges.

³ Note the comparison to other solvers using the propositional SAT formalism in Section 7.

Solution \mathcal{G}^1 s are represented via the predicate `edge1/2`, where `edge1(X,Y)` is *true* iff \mathcal{G}^1 contains the edge $X \rightarrow Y$. In ASP, the set of candidate solutions (i.e., the set of all directed graphs over n nodes) over which the search for solutions is performed, is declared via the so-called *choice construct* within the following rule, stating that candidate solutions may contain directed edges between any pair of nodes. If we have prior knowledge about edges that must (or must not) be present in \mathcal{G}^1 , then that content can straightforwardly be encoded here.

```
{ edge1(X,Y) } :- node(X), node(Y).
```

This is a so-called *choice rule* in the ASP syntax, which here states that `edge1` can be true or false for any pair of nodes X, Y , as given by the predicate `node`.

The implied measurement timescale structure \mathcal{G}^u for a candidate solution \mathcal{G}^1 is represented using the predicates `edgeu/2` and `confu/2`, which are derived in the following way. First, we declare the mapping from a given \mathcal{G}^1 to the corresponding \mathcal{G}^u by declaring the exact length- L paths in a non-deterministically chosen candidate solution \mathcal{G}^1 . For this, we declare rules that compute the length- L paths inductively for all $L \leq U$, using the predicate `path(X,Y,L)` to represent that there is a length- L path from X to Y .

```
% Derive all directed paths up to length U
path(X,Y,1) :- edge1(X,Y).
path(X,Y,L) :- path(X,Z,L-1), edge(Z,Y), L <= U, u(U).
```

The first rule states that an edge $X \rightarrow Y$ implies the existence of the (corresponding) path of length one. The second rule declares inductively, that the existence of a path of length $L - 1$ from X to Z , and an edge $Z \rightarrow Y$, together imply the existence of a path of length L from X to Y .

Second, to obtain \mathcal{G}^u , we encode Equations (1) and (2) with the following rules that form predicates `edgeu` and `confu` describing the edges \mathcal{G}^1 induces on the measurement timescale structure \mathcal{G}^u . The first rule derives induced directed edges in \mathcal{G}^u from the length- U paths, and the second the bidirected edges based on the existence of pairs of confounding paths of length up to $U - 1$.

```
% Paths of length U, correspond to measurement timescale edges
edgeu(X,Y) :- path(X,Y,U), u(U).

% Paths of equal length (<U) from a single node result in bi-directed edges
confu(X,Y) :- path(Z,X,L), path(Z,Y,L), node(X;Y;Z), X < Y, L < U, u(U).
```

Finally, we declare constraints that require that the \mathcal{G}^u represented by the `edgeu` and `confu` predicates is consistent with the input \mathcal{H} . This is achieved with the following *integrity* rules, which enforce that the edge relations of \mathcal{G}^u and \mathcal{H} are exactly the same for any solution \mathcal{G}^1 . In other words, the first two rules derive a contradiction in case the directed edge relations of \mathcal{G}^u and \mathcal{H} do not match; the third and fourth rules do the same for the bidirected edge relations of \mathcal{G}^u and \mathcal{H} . For example, if the `edgeh` is true in the input for some X and Y and the corresponding `edgeu` is not derived, the set of edges defined by `edge1` does not constitute a consistent graph for the input \mathcal{H} according to the first rule below.

```
:- edgeh(X,Y), not edgeu(X,Y).
:- no_edgeh(X,Y), edgeu(X,Y).
:- confh(X,Y), not confu(X,Y).
:- no_confh(X,Y), confu(X,Y).
```

Our ASP encoding of [Task 1](#) consists of the rules just described. The set of solutions of the encoding correspond exactly to the \mathcal{G}^1 s consistent with the input \mathcal{H} . Note that before solving, these first-order rules are grounded for all possible instantiations of X, Y, Z and L relevant to the input.

3.3. Runtime comparison

Both our proposed SAT-based approach and the recent specialized search algorithm MSL of Plis et al. [31] are correct and complete, so we focus on differences in efficiency, using the implementation of MSL by the original authors. Our approach

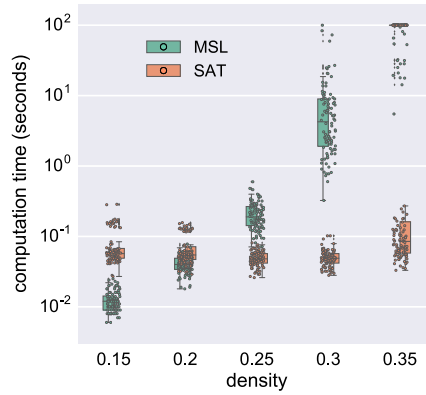


Fig. 3. Running times for 10-node rolled graphs as a function of graph density for the state of the art (MSL) and our method (SAT). We used 100 graphs per density and a timeout of 100 seconds; both methods enumerate up to 1000 solutions.

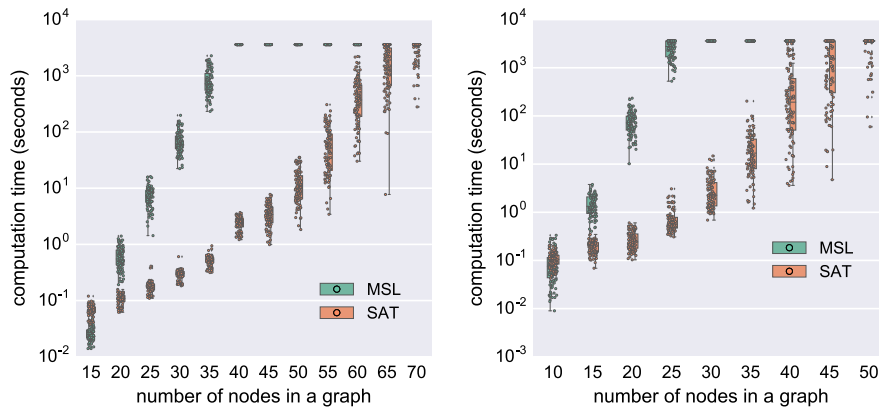


Fig. 4. Running times as function of the number of nodes for the state of the art (MSL) and our method (SAT). Left: 10%-dense graphs. Right: 15%-dense graphs. In both plots we use 100 graphs per size and a timeout of 1 hour; both methods enumerate up to 1000 solutions.

allows for searching simultaneously over a range of values of u , but Plis et al. [31] focused on the case $u = 2$; hence, we restrict the comparison to $u = 2$.

The MSL algorithm starts by noting that every measurement timescale edge corresponds to a path of length u in \mathcal{G}^1 , where that path must be through another measured variable. MSL thus creates $u - 1$ “virtual” mediating nodes for each measurement timescale edge, and then finds all ways of identifying virtual nodes with actual nodes such that all-and-only the measurement timescale edges are implied. Exhaustive search of all possible virtual to actual identifications is computationally intractable, so MSL employs a branch-and-bound search procedure, where a branch is bounded whenever it implies a “false positive” (i.e., implies an edge that does not actually occur in the measurement timescale input). Because each edge requires $u - 1$ virtual nodes, each of which must later be identified with an actual node, MSL scales quite poorly as a function of u .

For the comparison, we simulated system timescale rolled graphs with varying density and number of nodes (see Section 5 for exact details), and then computed the implied measurement timescale structures for subsampling rate $u = 2$. This structure was given as input to the inference procedures (including the subsampling rate $u = 2$). Note that the input consisted here of graphs for which there always is a \mathcal{G}^1 , so all instances were satisfiable. The task of the algorithms was to output up to 1000 (system timescale) graphs in the equivalence class. The ASP encoding was solved by Clingo using the flag `-n 1000` for the solver to enumerate 1000 solution graphs (or all, in cases where there were fewer than 1000 solutions).

The running times of the MSL algorithm and our approach (SAT) on 10-node (rolled) input graphs with different edge densities are shown in Fig. 3. Fig. 4 shows the scalability of the two approaches in terms of increasing number of nodes in the rolled input graphs and fixed 10% or 15% edge density. Our declarative approach clearly outperforms MSL. 10-node rolled input graphs, regardless of edge density, are essentially trivial for our approach, while the performance of MSL deteriorates noticeably as the density increases. For varying numbers of nodes in 10% density input graphs, our approach scales up to 65 nodes with a one hour time limit; even for 70 nodes, 25 graphs finished in one hour. In contrast, MSL reaches only 35 nodes; our approach uses only a few seconds for those graphs. The scalability of our algorithm allows for investigating the influence of edge density for larger graphs. Fig. 5 (left) plots the running times of our approach (when enumerating all

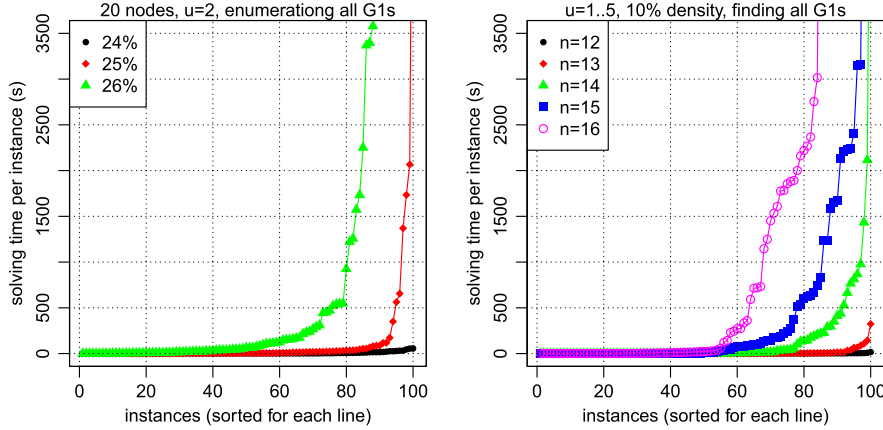


Fig. 5. Left: Influence of input graph density on running times of our approach when the subsampling rate $u = 2$ is given as input and all solutions are enumerated. Right: Scalability of our approach when u is left to be determined by the method from interval $1, \dots, 5$. All solutions over the range of u are enumerated.

solutions) for $u = 2$ ($u = 2$ was given as input) on 20-node input graphs of varying densities. Note that here the instances are sorted by the running time for each individual density (curve). With a time limit of 1000 seconds we can solve 80% of the instances with 26% density, almost all of the instances with 25% density and all of the instances with 24% density. Thus, the running time is increased for denser graphs: in addition to more constraints, there are also more members in the equivalence classes. Finally, Fig. 5 (right) shows the scalability of our approach in the more challenging task of enumerating *all* solutions over the range $u = 1, \dots, 5$ simultaneously. This also demonstrates the generality of our approach: it is not restricted to solving for individual values of u separately.

4. Learning system timescale structures from data

Due to statistical errors in estimating \mathcal{H} and the sparse distribution of implied \mathcal{G}^u in the space of possible undersampled graphs, the estimated \mathcal{H} will often have no \mathcal{G}^1 s with $\mathcal{G}^u = \mathcal{H}$. Given such an \mathcal{H} , neither the MSL algorithm nor our approach in the previous section can output a solution, and they simply conclude that no solution \mathcal{G}^1 exists for the input \mathcal{H} .⁴ In terms of our constraint declarations, this is witnessed by conflicts among the constraints and the underlying model space for any possible solution candidate. Given the inevitability of statistical errors, we should not simply conclude that no consistent \mathcal{G}^1 exists for such an \mathcal{H} . Rather, we should aim to learn \mathcal{G}^1 s that, in light of the underlying conflicts, are “optimally close” (in some well-defined sense of “optimal”) to being consistent with \mathcal{H} . We now turn to this more general problem setting, and propose what (to the best of our knowledge) is the first approach to learning, by employing constraint optimization, from undersampled data under conflicts. In fact, we can use the ASP formulation already discussed—with minor modifications—to address this problem.

In this more general setting, the input consists of both the estimated graph \mathcal{H} , and also (i) weights $w(e \in \mathcal{H})$ indicating the reliability of edges present in \mathcal{H} ; and (ii) weights $w(e \notin \mathcal{H})$ indicating the reliability of edges absent in \mathcal{H} . Since \mathcal{G}^u is \mathcal{G}^1 subsampled by u , the task is to find a \mathcal{G}^1 that minimizes the objective function

$$f(\mathcal{G}^1, u) = \sum_{e \in \mathcal{H}} I[e \notin \mathcal{G}^u] \cdot w(e \in \mathcal{H}) + \sum_{e \notin \mathcal{H}} I[e \in \mathcal{G}^u] \cdot w(e \notin \mathcal{H}),$$

where the indicator function $I(c) = 1$ if the condition c holds, and $I(c) = 0$ otherwise. Thus, edges that differ between the estimated input \mathcal{H} and the \mathcal{G}^u corresponding to the solution \mathcal{G}^1 are penalized by the weights representing the reliability of the measurement timescale estimates. In the following, we first outline how to generalize the ASP encoding from the preceding section to enable search for optimal \mathcal{G}^1 with respect to this objective function. We then describe two alternatives for determining the weights w . In the following section, we present simulation results on the relative performance of the different weighting schemes.

4.1. Learning by constraint optimization

To model the objective function for handling conflicts, only simple modifications are needed to our ASP encoding: instead of declaring *hard* constraints that require that the paths induced by \mathcal{G}^1 *exactly* correspond to the edges in \mathcal{H} , we *soften*

⁴ For these cases, Plis et al. [31] ran MSL on graphs close to \mathcal{H} to try to find an input for which there is a \mathcal{G}^1 , but this strategy is not guaranteed to find an optimal solution, nor does it scale computationally.

these constraints by declaring that the violation of each individual constraint incurs the associated weight as penalty. In the ASP language, this can be expressed by augmenting the input predicates $\text{edge}_{\mathcal{H}}(X, Y)$ with weights: $\text{edge}_{\mathcal{H}}(X, Y, W)$ (and similarly for $\text{no_edge}_{\mathcal{H}}$, $\text{conf}_{\mathcal{H}}$ and $\text{no_conf}_{\mathcal{H}}$), and by using *weighted soft rules* syntactically represented via $:\sim$ instead of $:-$. Here the additional argument W represents the weight $w((X \rightarrow Y) \in \mathcal{H})$ given as input. The following expresses that each conflicting presence of an edge in \mathcal{H} and \mathcal{G}^u is penalized with the associated weight W . The additional $[W, X, Y, v]$ for $v = 1, 2$ syntactically enforce that a cost of W is incurred in case the corresponding rule is violated for a specific pair of nodes X, Y . The numbers $v \in \{1, 2\}$ at the end of the brackets enable the solver to distinguish the cost incurred due to bidirected and directed edges respectively.

```

: $\sim$  edge $_{\mathcal{H}}(X, Y, W)$ , not edge $_{\mathcal{G}^u}(X, Y)$ . [W, X, Y, 1]
: $\sim$  no_edge $_{\mathcal{H}}(X, Y, W)$ , edge $_{\mathcal{G}^u}(X, Y)$ . [W, X, Y, 1]
: $\sim$  conf $_{\mathcal{H}}(X, Y, W)$ , not conf $_{\mathcal{G}^u}(X, Y)$ . [W, X, Y, 2]
: $\sim$  no_conf $_{\mathcal{H}}(X, Y, W)$ , conf $_{\mathcal{G}^u}(X, Y)$ . [W, X, Y, 2]

```

This modification provides an ASP encoding for [Task 2](#); that is, the optimal solutions to this ASP encoding correspond exactly to the \mathcal{G}^1 s that minimize the objective function $f(\mathcal{G}^1, u)$ for given u and input \mathcal{H} with weighted edges.

4.2. Weighting schemes

We use two different schemes for weighting the presences and absences of edges in \mathcal{H} according to their reliability. To determine the presence or absence of a specific edge $X \rightarrow Y$ in \mathcal{H} , we simply test the corresponding independence $X^{t-1} \perp\!\!\!\perp Y^t \mid \mathbf{V}^{t-1} \setminus X^{t-1}$. To determine the presence/absence of an edge $X \leftrightarrow Y$ in \mathcal{H} , we test the independence: $X^t \perp\!\!\!\perp Y^t \mid \mathbf{V}^{t-1}$.

The simplest approach is to use uniform weights for the estimated \mathcal{H} :

$$w(e \in \mathcal{H}) = 1 \quad \forall e \in \mathcal{H},$$

$$w(e \notin \mathcal{H}) = 1 \quad \forall e \notin \mathcal{H}.$$

Uniform edge weights resemble the search on the Hamming cube of \mathcal{H} that Plis et al. [31] used to address the problem of finding \mathcal{G}^1 s when \mathcal{H} did not correspond to any \mathcal{G}^u , though our approach is much superior computationally.

A more intricate approach is to use pseudo-Bayesian weights following [24,16,34]. They used Bayesian model selection to obtain reliability weights for independence tests. Instead of a p -value and a binary decision, these types of tests give a measurement of reliability for an independence/dependence statement as a Bayesian probability. We can directly incorporate their approach of using log-probabilities as the reliability weights for the edges. For details, see Section 4.3 of Hyttinen et al. [16]. Again, we only compute weights for the independence tests mentioned above in the estimation of \mathcal{H} .

5. Simulations

We use simulations to explore the accuracy and runtime efficiency of our approach in various different settings. For the simulations, system timescale structures \mathcal{G}^1 and the associated data generating models were constructed in the following way. To guarantee connectedness of the graphs, we first formed a cycle of all nodes in a random order (following Plis et al. [31]). We then randomly sampled additional directed edges until the required density was obtained. Recall that there are no bidirected edges in \mathcal{G}^1 . We used Equations (1) and (2) to generate the measurement timescale structure \mathcal{G}^u for a given u . When sample data were required, we used linear Gaussian structural autoregressive processes (order 1) with structure \mathcal{G}^1 to generate data at the system timescale, where coefficients were sampled from the two intervals $\pm[0.2, 0.8]$. We then discarded intermediate samples⁵ to get the particular subsampling rate.⁶

5.1. Accuracy

[Fig. 6](#) shows the accuracy of the different methods in one setting: subsampling rate $u = 2$ (given as input), network size $n = 6$, average degree 3 (density 25%), $N = 250$ samples, and 200 datasets in total. The positive predictions correspond to presences of edges; when the method returned several solutions with equal cost, we used the mean solution accuracy to measure the output accuracy. The x-axis numbers correspond to the adjustment parameter for the statistical independence tests (prior probability of independence). The two left columns (black and red) show the true positive rate and false positive rate of the \mathcal{H} estimation (compared to the true \mathcal{G}^2), for the different types of edges, using different statistical tests. Given 250 samples, we see that the structure of \mathcal{G}^2 can be estimated with a good tradeoff of TPR and FPR with the middle

⁵ All sample counts refer to the number of samples after subsampling.

⁶ `Clingo` only accepts integer weights; we multiplied weights by 1000 and rounded to the nearest integer.

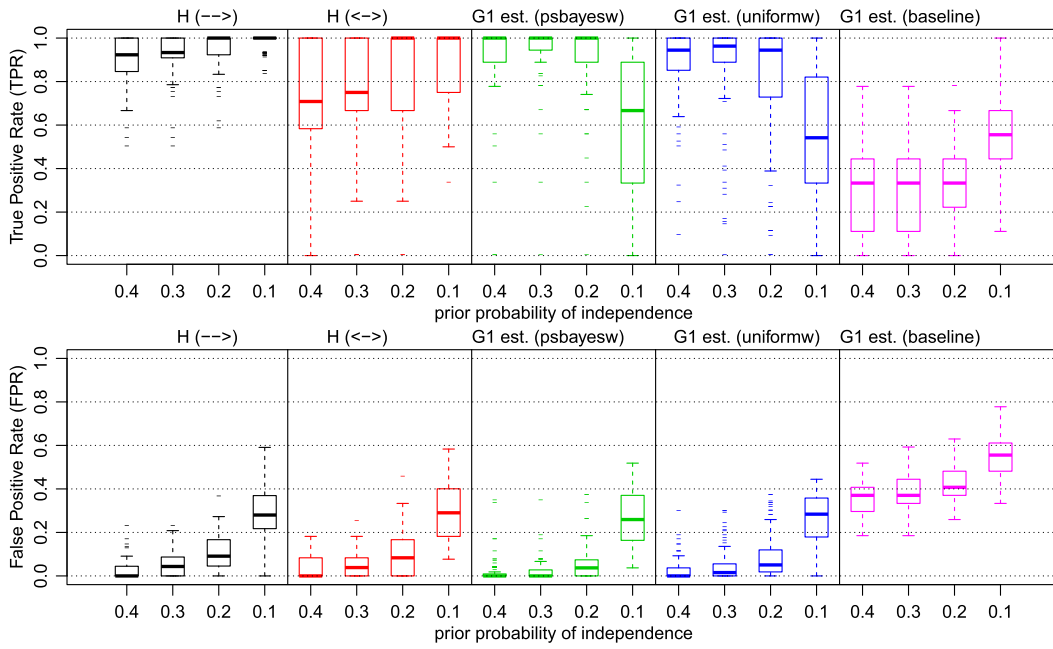


Fig. 6. Accuracy of the optimal solutions when subsampling rate $u = 2$ is given as input (200 instances and 250 samples). The x-axis shows the different prior probabilities of independence in the utilized independence test. The two left columns give the accuracy of the estimation of the measurement timescale structure \mathcal{H} . The next two columns give the accuracy of our method with the two different weighting schemes. The rightmost column shows the accuracy of the baseline estimate that does not take subsampling into account (the directed edges of \mathcal{H} are directly interpreted as the system timescale edges). (For interpretation of the references to color in this figure, the reader is referred to the web version of this article.)

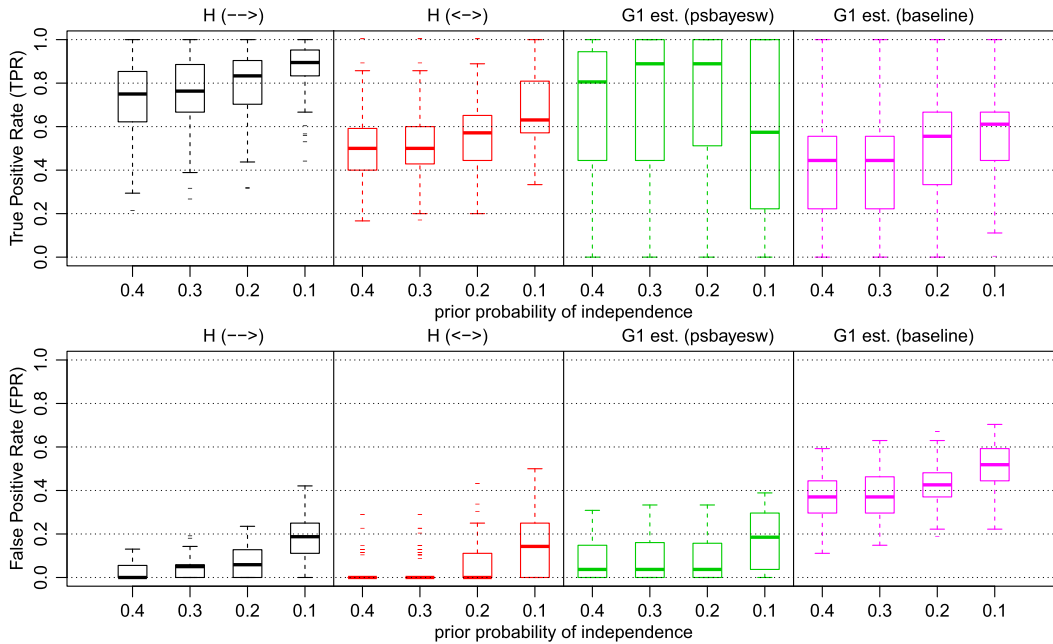


Fig. 7. Accuracy of the optimal solutions when subsampling rate $u = 2$ is given as input (200 instances and 500 samples). The x-axis shows the different prior probabilities of independence in the utilized independence test. The two left columns give the accuracy of the estimation of the measurement timescale structure \mathcal{H} . The third column gives the accuracy of our method with the pseudo-Bayesian weighting scheme. The rightmost column shows the accuracy of the baseline estimate that does not take subsampling into account. (For interpretation of the references to color in this figure, the reader is referred to the web version of this article.)

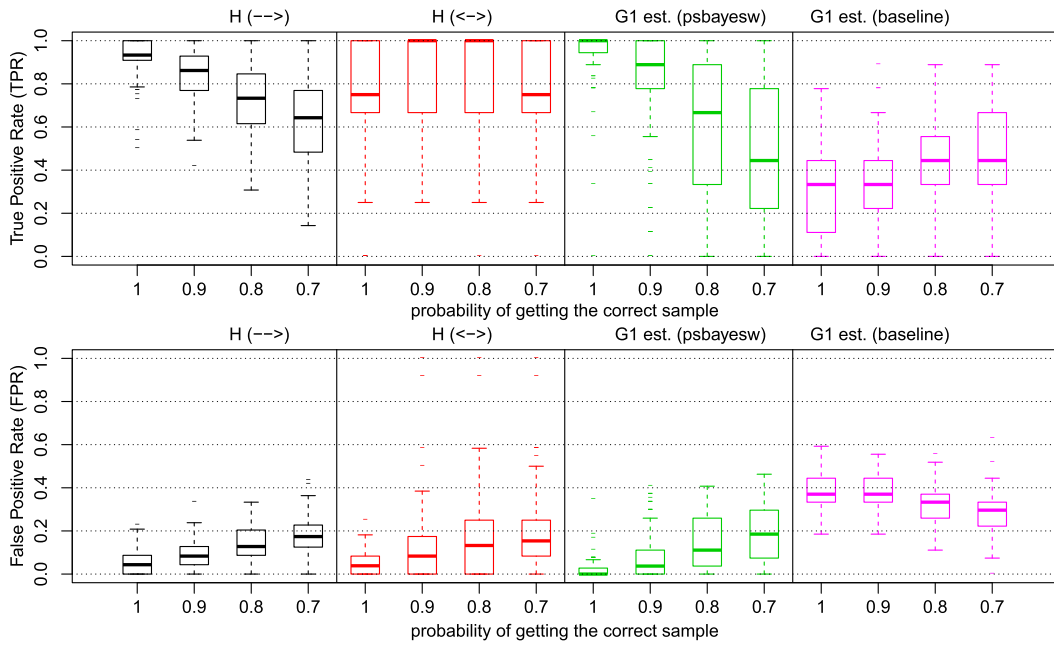


Fig. 8. Accuracy of the optimal solutions when subsampling rate $u = 2$ is given as input (200 instances and 250 samples), some samples are obtained at the adjacent timepoints. Due to previous simulations we used the prior probability of 0.3 for all methods. In more detail, the x-axis gives the probability that the sample was obtained at the correct time t , otherwise the sample was obtained either at the previous or the next time point, splitting the remaining probability. The two left columns give the accuracy of the estimation of the measurement timescale structure \mathcal{H} . The third column gives the accuracy of our method with the pseudo-Bayesian weighting scheme. The rightmost column shows the accuracy of the baseline estimate that does not take subsampling into account.

parameter values, but not perfectly. The presence of directed edges can be estimated more accurately. More importantly, the third and fourth columns in Fig. 6 (green and blue) show the accuracy of the \mathcal{G}^1 estimation. Both weighting schemes produce good accuracy for the middle parameter values, although there are some outliers. The pseudo-Bayesian weighting scheme (“psbayesw”, shown in green) still outperforms the uniform weighting scheme (“uniformw”, shown in blue), as it produces high TPR with low FPR for a range of threshold parameter values (especially for 0.3). Both weighting schemes are superior to the “baseline” shown in magenta on the right. This baseline \mathcal{G}^1 estimate is formed by the directed edges of the estimated H , and thus corresponds to estimating \mathcal{G}^1 without taking subsampling into account.

Fig. 7 shows the accuracy when $u = 3$ (given as input), $n = 6$, average degree 3 (density 25%), $N = 500$, and 200 datasets. The accuracy for edge presences in the measurement timescale graph \mathcal{H} is lower than for $u = 2$, even though we have twice the number of samples (Fig. 7, two rightmost columns in black and red). The problem is that measurement timescale edges here correspond to 3-edge paths, whose causal effects will be smaller (on average) than 2-edge paths for a fixed interval of system timescale edge coefficients ($\pm[0.2, 0.8]$), and so are harder to detect. Nevertheless, the constraint optimization procedure achieves a good tradeoff between TPR and FPR for system timescale edges (Fig. 7, third column in green). Larger subsampling rates (u) require more samples for accurate \mathcal{G}^1 structure discovery, but not several orders of magnitude more data.

5.2. Robustness of the subsampling rate

Fig. 8 shows the accuracy of this method when some of the samples are not obtained at the exact time assumed by the measurement timescale. Specifically, the x-axis specifies the probability with which we obtain the correct sample (for the given $u = 2$, which is given as input); otherwise, we take either the sample before or the sample after (synchronously for all variables), splitting the remaining probability. The results with probability 1 equal the result in Fig. 6 with prior probability of independence 0.3 and $N = 250$ samples. These values were used in all runs in this plot. Unsurprisingly, as the “jitter” in the sampling process increases, the results deteriorate in terms of TPR and FPR. However, at least for the models and subsampling rate of $u = 2$ tested here, the inference is not overly sensitive. When the probability of a correct sample is 0.9, the results are still quite good, alleviating somewhat the dependence on the assumption of an exact subsampling rate. Naturally, there are many further permutations one could explore: jitter could affect variables independently of one another, jitter could be represented by a more complex distribution, we could explore the effect of jitter for different subsampling rates or when the subsampling rate is unknown. Moreover, jitter could have a persistent, rather than a local effect, in shifting subsequent measures as well. We have here only explored the simple case mimicking the situation where the measurement device as a whole (i.e. simultaneously for all variables) comes out of synch with the system at random points without consequences for subsequent samples.

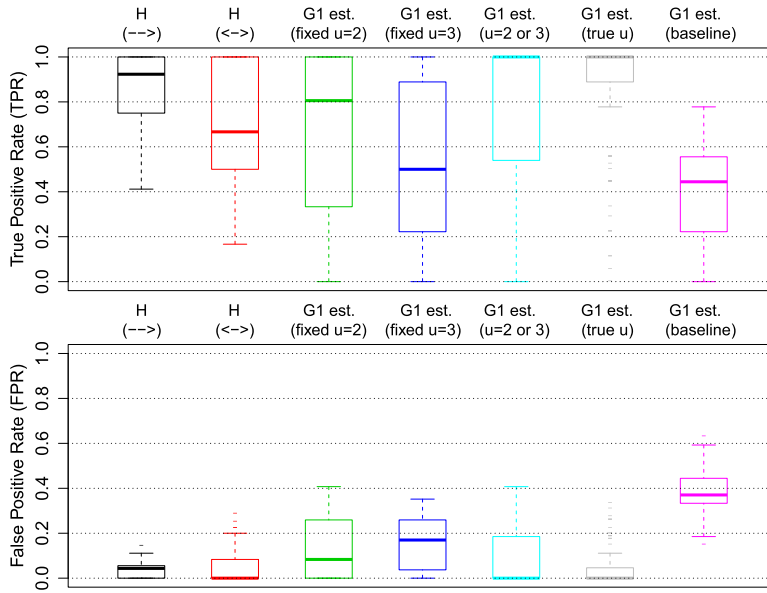


Fig. 9. Accuracy when the true u is unknown. Two left boxplots show accuracy of the \mathcal{H} estimate as before. The next three boxplots show the accuracy of our approach (pseudo-Bayesian weights) when, regardless of the true u , u is fixed to 2, or to 3, or left for the procedure decision, respectively. In the second from right boxplot the true u was given as input, the rightmost boxplot shows the baseline that does not take subsampling into account.

Fig. 9 further examines the possibility to distinguish between different subsampling rates. We generated 500 samples of data from 200 models (average degree 3) with equal numbers of cases with $u = 2$ or $u = 3$. The two leftmost boxplots show the accuracy of the estimated \mathcal{H} , which, given the mixture of $u = 2$ and $u = 3$, is between the accuracy of \mathcal{H} obtained in previous simulations. The next two boxplots show the accuracy of the \mathcal{G}^1 estimate, when the subsampling rate u for the search procedure is fixed to 2 or 3, respectively, regardless of the true u . As expected, the accuracy is mediocre in this case, since the method assumes the incorrect subsampling rate u in half of the runs. But when the method is left to determine the correct u by itself, the accuracy improves again, as shown in the boxplots third from right (the method was run with $u = 2 \dots 3$). In fact, the accuracy comes close to that of the second from right boxplot, where the correct u was given as input to the procedure. Thus the procedure is often able to recognize the correct u . The longer tails indicate that at times the determination of u is not perfect.

5.3. Scalability

Finally, the running times of our approach are shown in Fig. 10 with different weighting schemes, network sizes (n), and numbers of samples (N). The subsampling rate was again fixed to $u = 2$ (and given as input), and average node degree was 3. Fig. 10 (left) shows that the pseudo-Bayesian weighting scheme allows for much faster solving: for $n = 7$, it finishes all runs in a few seconds (black circle), while the uniform weighting scheme (red diamond) takes several minutes in the longest runs. Thus, the pseudo-Bayesian weighting scheme provides the best performance in terms of both computational efficiency and accuracy. The number of samples has a significant effect on the running times: larger number of samples take less time. Runs for $n = 9$, $N = 200$ (blue square) take longer than for $n = 9$, $N = 500$ (Fig. 10 left, magenta circle vs. cyan diamond). Intuitively, statistical tests should be more accurate with larger number of samples, resulting in fewer conflicting constraints. For $N = 1000$, the global optimum is found here for up to 12-node graphs (Fig. 10 right), though in a considerable amount of time.

6. Case study: house data of Peters et al. [29]

In order to demonstrate the applicability to real-world data, we analyzed the house temperature and humidity data of Peters et al. [29]. The data includes 7265 samples of hourly temperature and humidity measurements of six sensors placed in a house (SHED = in the shed, OUT = outside, KIT = kitchen boiler, LIV = living room, WC = wc, BATH = bathroom) in the Black Forest. The house has heating, but the house is not in use for most of the year. This data was also partly analyzed by Gong et al. [11]. The measurements of this system were obtained at coarser intervals than the process of temperature and humidity changes are thought to take place. Since the data includes outside temperature and humidity measurements, the assumption of causal sufficiency at the system timescale seems a good approximation.

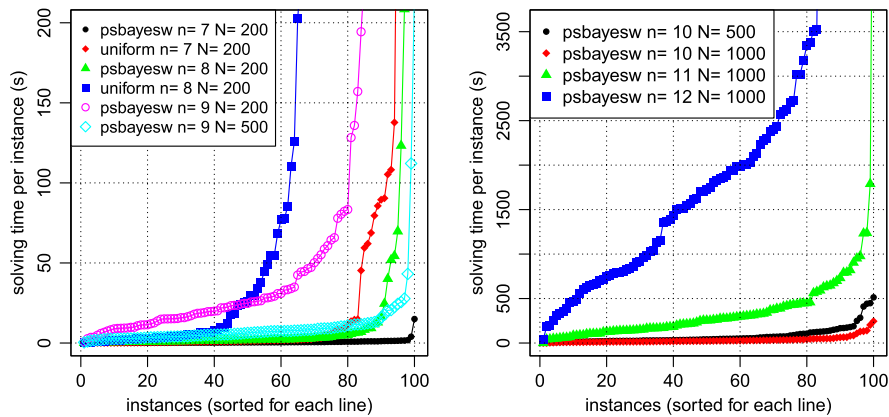


Fig. 10. Scalability of our constraint optimization approach (using `Clingo`) for different graph sizes, numbers of samples and weighting schemes. For each setting there are 100 instances that are sorted according to the solving time on each line. (For interpretation of the references to color in this figure, the reader is referred to the web version of this article.)

We analyzed the temperature and humidity components separately, and examined the differences of sequential measurements,⁷ as this removed trends from each univariate time series. The temperature measurement timescale graph (obtained at 0.9 prior probability of independence) includes a total of 20 (out of 36) directed edges, and 8 (out of 15) bidirected edges, with varying pseudo-Bayesian weights. The humidity measurement timescale graph had the same total numbers of edges, although not the exact same edges.

As explained earlier, subsampling introduces underdetermination of the system timescale graph. Thus, we determined the presence of individual system timescale edges in the following way [23]. For each edge in \mathcal{G}^1 , we ran the inference procedure first enforcing its presence and then enforcing its absence.⁸ The difference in objective function values for the two outputs—the optimal \mathcal{G}^1 s that do or do not contain the edge, respectively—indicates the support for the presence (absence) of the edge.

For the estimated \mathcal{H} , we computed \mathcal{G}^1 s edgewise for subsampling rates of $u = 2, 3$. (Since the measurements were hourly, these correspond to time steps of 30 and 20 minutes, respectively.) The two temperature graphs for $u = 2$ and $u = 3$ (Fig. 11a, b) differ substantially from one another, as do the two humidity graphs (Fig. 11d, e). These results provide empirical demonstrations of the impact of subsampling, as different choices of u imply different structures. At the same time, timesteps of 20 and 30 minutes arguably do not correspond to realistic time steps for the temperature and humidity changes measured by these data.

We thus considered larger subsampling rates $u = 10, 12$, which correspond to more realistic time steps of 5–6 minutes. As expected, there is more underdetermination for these u , but the results are also more plausible. Fig. 11c suggests that the temperature outside is not directly influenced by the temperature in any of the rooms, but it directly influences the temperature in the shed. The data do not, however, uniquely determine how the outside temperature directly affects the temperatures in the rooms inside the house, nor the system timescale causal dependencies between temperatures in the rooms. The algorithm output is both intuitively sensible, and also points towards future targeted experiments if the remaining underdetermination is to be resolved.

Similarly, the humidity structures for larger u are more plausible. Fig. 11f suggests that the humidity level in the WC is driven by both bathroom and outside humidity, which is sensible since the WC is located next to the bathroom and has a window, according to Peters et al. [29]. Similarly as Peters et al. [29], we find that the shed humidity affects bathroom humidity—for both analyses this may be due to an inability to distinguish the shed humidity from the outside humidity (they are particularly strongly correlated). The living room and kitchen boiler humidities seem to depend on each other directly, so the data suggest that the rooms may be adjacent, though that information was not provided by Peters et al. [29]. The algorithm thus points to testable predictions about the spatial house layout, and the mechanisms for humidity transfer.

Overall, the processes controlling the temperature and humidity have differences and similarities. Determining the placement of sensors thus seems to require data from both measurement types. More importantly for our present paper, this case study shows that this algorithm can be applied to real-world data, provide intuitively sensible outputs, and provide novel experiments and measurements that would resolve remaining underdetermination.

⁷ This may take out some of the influences of self-loops.

⁸ This can be done by adding a simple clause to the input code “`edge(X, Y) .`” to enforce the presence and “`:-edge(X, Y) .`” to enforce the absence of $X \rightarrow Y$.

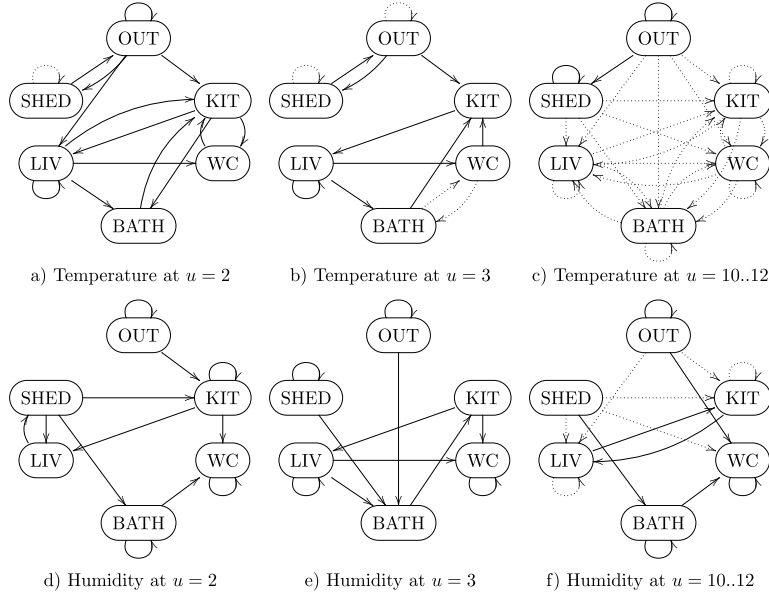


Fig. 11. Results of the House data analysis for different subsampling rates (u) and measurement type. Edges with full lines are found to be present, absent edges are found to be absent, edges with dotted lines may be present or absent.

7. Solver performance comparison

Thus far in this article we have considered `Clingo` as the only solver to find solutions to a declarative constraint encoding of the computational problems considered here. This raises the question to what extent the choice of the constraint solver affects the runtime performance of our approach. While the high-level ASP syntax is relatively easy to understand and modify, our approach can also be represented via propositional logic. The benefit of using propositional logic is that various SAT solvers, as well as MaxSAT solvers (as the Boolean optimization generalization of SAT), can be applied directly. In this section we evaluate the impact of the choice of SAT and MaxSAT solvers on the runtime efficiency of our approach.

7.1. Direct propositional SAT encoding

A direct propositional SAT encoding for finding a system timescale causal structure \mathcal{G}^1 consistent with a measurement timescale graph \mathcal{H} for a known u is presented in Eqs. (5)–(12).

$$\vec{h}_{X,Y} \quad \forall X, Y \in \mathbf{V} : X \rightarrow Y \in \mathcal{H} \quad (5)$$

$$\neg \vec{h}_{X,Y} \quad \forall X, Y \in \mathbf{V} : X \rightarrow Y \notin \mathcal{H} \quad (6)$$

$$\leftrightarrow h_{X,Y} \quad \forall X, Y \in \mathbf{V} : X < Y, X \leftrightarrow Y \in \mathcal{H} \quad (7)$$

$$\neg \leftrightarrow h_{X,Y} \quad \forall X, Y \in \mathbf{V} : X < Y, X \leftrightarrow Y \notin \mathcal{H} \quad (8)$$

$$\vec{h}_{X,Y} \Leftrightarrow \bigvee_{Z \in \mathbf{V}} (p_{X,Z}^{u-1} \wedge p_{Z,Y}^1) \quad \forall X, Y \in \mathbf{V} \quad (9)$$

$$p_{X,Y}^{l+1} \Leftrightarrow \bigvee_{Z \in \mathbf{V}} (p_{X,Z}^l \wedge p_{Z,Y}^1) \quad \forall X, Y \in \mathbf{V}, l \in \{1..u-2\} \quad (10)$$

$$\leftrightarrow h_{X,Y} \Leftrightarrow \bigvee_{l=1}^{u-1} \leftrightarrow h_{X,Y}^l \quad \forall X, Y \in \mathbf{V} : X < Y \quad (11)$$

$$\leftrightarrow h_{X,Y}^l \Leftrightarrow \bigvee_{Z \in \mathbf{V}} (p_{Z,X}^l \wedge p_{Z,Y}^1) \quad \forall X, Y \in \mathbf{V} : X < Y, l \in \{1..u-1\} \quad (12)$$

Essentially, Eqs. (5)–(8) enforce the input constraints imposed by \mathcal{H} . Following the ASP encoding presented earlier, Eqs. (9)–(12) encode the mapping from the \mathcal{G}^1 's—the edge relation of which is encoded as the length-1-path variables $p_{X,Y}^1$ —that are consistent with \mathcal{H} .

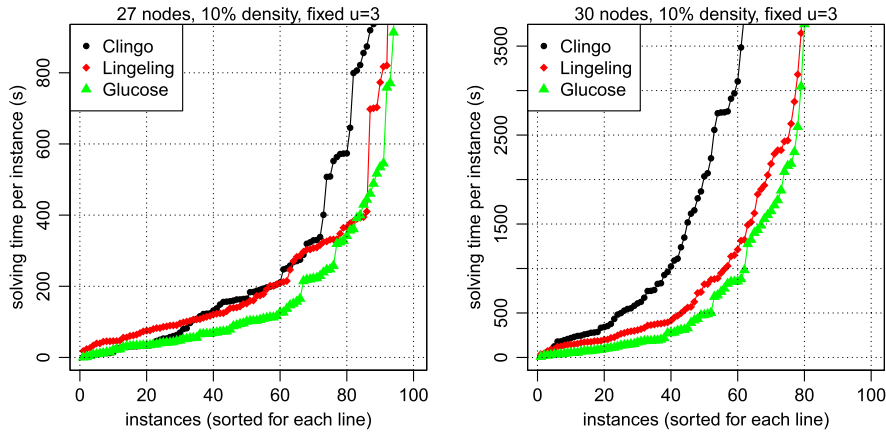


Fig. 12. Comparison of running times for different solvers finding a single graph in the equivalence class, when the subsampling rate $u = 3$ is given as input. Left: easier instances with 27 nodes. Right: harder instances with 30 nodes. `Clingo` uses the ASP encoding presented in Section 3.2, all others use the propositional SAT encoding in Section 7.1.

7.2. Solver comparison: finding consistent system timescale structures

The results of a runtime performance comparison between `Clingo` and two state-of-the-art SAT solvers, Glucose [2] and Lingeling [3], is presented in Fig. 12 for $u = 3$ (given as input), edge density of 10% and the numbers of nodes ranging from 27 (on left) to 30 (on right). Note that the plots give the running times of each of the three solvers sorted individually for each solver. In terms of runtime performance, the SAT solvers Glucose and Lingeling, both working directly on the propositional SAT encoding, exhibit noticeably improved performance over `Clingo` as the number of nodes is increased (right plot). Thus, in terms of runtime efficiency of our approach, it can be beneficial to apply current and future advances in state-of-the-art SAT solvers directly on the propositional level for improved performance. In these simulations the ASP paradigm does not show any particular computational advantage.

7.3. Solver comparison: learning system timescale structures from data

As with the ASP encoding given earlier, the SAT encoding given as Eqs. (5)–(12) is easily extended to solve the optimization problem underlying the task of learning system timescale structure from undersampled data. In the language of MaxSAT, the only change required is to make the constraints in Eqs. (5)–(8) soft, and to declare that the cost incurred from not satisfying these individual constraints equals that of $w(e \in \mathcal{H})$ (for Eqs. (5), (7)) or $w(e \notin \mathcal{H})$ (for Eqs. (6), (8)) for the corresponding edge e . This enables a comparison of the runtime performance of `Clingo`'s default branch-and-bound based search for an optimal solution to those of other MaxSAT solvers implementing alternative algorithmic approaches on the direct propositional MaxSAT encoding. Results comparing the performance of `Clingo` to that of the modern MaxSAT solvers Eva500a [27], LMHS [32], MSCG [26], Open-WBO [25], PrimalDual [5], and QMaxSAT [20], as well as the commercial integer programming (IP) solver CPLEX run on a standard IP translation of MaxSAT [8,1], are shown in Fig. 13. Here we observe that `Clingo`'s branch-and-bound approach is among the best performing solvers (with the considered problem parameters). However, the results also suggest that QMaxSAT, and so-called model-based approaches using a SAT solver to search for an optimal solution over the objective function range with a top-down strategy, can improve on the runtime efficiency of our approach. These results clearly show that the choice of the underlying Boolean optimization solver can indeed have a noticeable influence on the practical efficiency of the approach. There is at least some potential for further improving the runtime performance of our approach by making use of advances in MaxSAT solver technology.

8. Learning from mixed frequency data

In some contexts we may have obtained data from the same system at different subsampling frequencies. Two cases can be distinguished here: First, the subsampled time series may be anchored to the same underlying process such that one may know about the offset between the two.⁹ For approaches to this case see Tank et al. [36], who treat this issue as a missing data problem in a parametric setting. The second case we consider here is one where the subsampled time series are taken at different times and cannot be coordinated to the same instance of an underlying time series. A natural question is how much more can be learned by integrating information from multiple sampling rates. If one sampling rate is an integer multiple of the other, then (provably) nothing additional can be learned. A more interesting situation arises when

⁹ For example, in the special case with two simultaneously measured data sets with $u = 2$ and 1 time step offset, we can combine the time series to give a dataset with no subsampling.

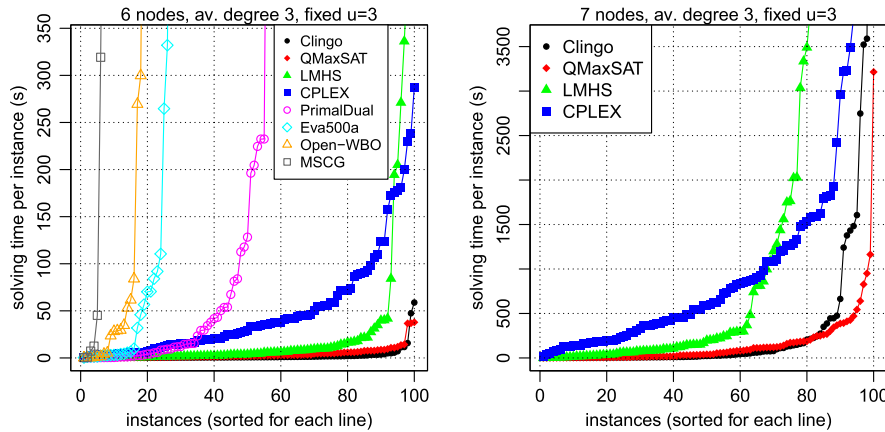


Fig. 13. Comparison of running times for different solvers finding the optimal graph, when the subsampling rate $u = 3$ is given as input. Left: easier instances with 6 nodes. Right: harder instances with 7 nodes. Clingo uses the ASP encoding presented in Sections 3.2 and 4.1, all others use the propositional SAT encoding in Section 7.1.

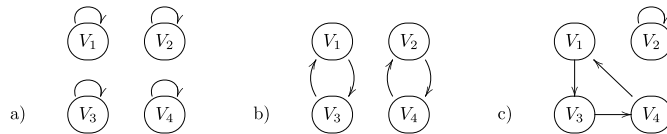


Fig. 14. Example graphs for learning from mixed frequency data. Graph (a) shows the true system timescale causal structure. When this is subsampled by $u = 2$ or by $u = 3$, the result is also the structure (a) (this time in measurement timescale). System timescale structure (b) gives measurement timescale structure (a) when subsampling by $u = 2$. System timescale structure (c) gives measurement timescale structure (a) when subsampling by $u = 3$. However, if measurement timescale structures for $u = 2$ and $u = 3$ are given as (a) respectively, the true system timescale structure can in fact be identified as (a).

neither sampling rate is an integer multiple of the other. For example, suppose the causal system operates at a 1-second timescale. If the system is measured every 2 seconds in one dataset, and every 3 seconds in another dataset, then we have $u_1 = 2/3 \cdot u_2$. More generally, if u_1/u_2 is non-integer, then when (if ever) is the equivalence class of \mathcal{G}^1 that satisfies both \mathcal{H}_1 & \mathcal{H}_2 smaller than the equivalence class for either \mathcal{H} individually? We can start to answer this question using the constraint satisfaction approach of this paper with only minor modifications.

For example, suppose the true system timescale structure is given in Fig. 14a. That is, the system includes four independent time series with self-loops. Undersampling does not change this graph, so the measurement timescale structures for $u = 2$ and for $u = 3$ will also be the graph in Fig. 14a. For this measurement timescale graph, the system timescale structure is not uniquely determined for either $u = 2$ or $u = 3$: for example, the system timescale structure in Fig. 14b produces Fig. 14a with $u = 2$, and Fig. 14c produces Fig. 14a with $u = 3$. In fact, any system timescale edge can be present or absent given either of the measurement timescale graphs alone.¹⁰ However, if this measurement timescale graph is found at both $u = 2$ and $u = 3$, then the system timescale structure can be uniquely determined: Fig. 14b produces a different measurement timescale graph for $u = 3$ and Fig. 14c produces a different measurement timescale graph for $u = 2$. And of course, the same observations hold if the u s are multiplied by a constant (e.g., if $u = 4$ and $u = 6$).

To examine the prevalence of this phenomenon, we exhaustively considered all $65536 (= 2^{4^4})$ different 4-variable \mathcal{G}^1 s, and compared the number of equivalence classes given input at a single subsampling rate, versus given inputs at two subsampling rates. A greater number of equivalence classes means a higher chance that a random graph will be uniquely identifiable, and so the number of equivalence classes is an approximate (inverse) measure of the extent of underdetermination.

For input at a single undersampling rate, for $u = 2$ we have 24265 equivalence classes; 7544 for $u = 3$; and 3964 equivalence classes for $u = 4$. These results with a single undersampled input graph thus replicate the known result that underdetermination is a significant problem, and it rapidly worsens as u increases [30,31].

If we instead have measurement timescale graphs for both $u = 2, 3$, then we have 26720 equivalence classes, which is only slightly more than the number for $u = 2$ by itself. That is, underdetermination is not substantially reduced if we additionally measure at $u = 3$ when we already have measurements at $u = 2$. Similarly, for $u = 3, 4$ we have 7814 equivalence classes; again, there is a reduction in underdetermination compared to $u = 3$ by itself, but it is quite small. This analysis assumes that all \mathcal{G}^1 are equally likely, and it is an open question whether measurements at different undersampling rates would have more impact for certain classes of \mathcal{G}^1 (e.g., connected graphs).

¹⁰ The node labels in Fig. 14b and c can be permuted.

9. Discussion

We have assumed that all common causes of measured variables are themselves measured, but this assumption is frequently violated in real-world data. Constraint satisfaction methods have elsewhere been used with success to identify causal relations in the presence of unobserved common causes or latent variables [16,23]. For time series data, dropping the assumption of causal sufficiency (in the system timescale) generates complications. Even if the system timescale process including latent variables is assumed to be first order Markov, the Markov order of the measurement timescale (naturally without the latent variables) can be arbitrarily larger.¹¹ That is, variables arbitrarily far in the past can (directly, in the measurement timescale) cause variables at the current timestep. We would thus need to both enrich the notation for \mathcal{G}^u to encode the time lags of direct causal effects, and also modify the statistical tests used to estimate these connections.

Moreover, there can be more information contained in the pattern of time lags (i.e., which past variables directly cause the present) than is given by the Markov order of the system. As just one example, suppose $\{X^{t-2}, X^{t-4}, \dots\} \rightarrow Y^t$. The simplest (in terms of number of latents) structure that explains these influences (i) has a latent L through which X influences Y (i.e., $X^{t-2} \rightarrow L^{t-1} \rightarrow Y^t$); and (ii) L is part of a 2-loop with another latent M (i.e., $L^{t-1} \rightarrow M^t$ and $L^t \leftarrow M^{t-1}$). In contrast, if we have $\{X^{t-2}, X^{t-3}, \dots\} \rightarrow Y^t$, then the simplest structure has only a single latent L through which X influences Y , but where L has a self-loop (i.e., $L^{t-1} \rightarrow L^t$). The pattern of time lags for direct causes—in particular, the absence of certain time lags—thus contains information about the number and causal structure of the latent variables. Estimation of this pattern, however, can be quite complex statistically.

Subsampled time series data can be also particularly prone to violations of faithfulness. For example, the underlying process unrolled over time may include directed paths over many time steps that do not result in significant statistical dependence in the observed data. In addition, variables observed over subsequent time steps might be almost deterministically related. If $X^{t-1} \approx X^{t-2}$, then conditioning on X^{t-2} may render the statistical dependence through $Y^t \leftarrow X^{t-1} \rightarrow Z^t$ undetectable from any realistic numbers of samples. In the current framework, both of these situations are treated as estimation errors in \mathcal{H} . Further modeling of these complications may help to achieve improved accuracy. Another option could be to develop parametric approaches instead of the non-parametric one presented in this paper.

10. Conclusion

In this paper, we introduced a constraint optimization based solution for the problem of learning causal timescale structures from subsampled measurement timescale graphs and data. Our approach considerably improves the state-of-art; in the simplest case (subsampling rate $u = 2$), we extended the scalability by several orders of magnitude. Moreover, our method generalizes to handle different or unknown subsampling rates in a computationally efficient manner. Unlike previous methods, our method can operate directly on finite sample input, and we presented approaches that recover, in an optimal way, from conflicts arising from statistical errors. We demonstrated the accuracy, robustness and scalability of the approach through a series of simulations and applied it to real-world time series data. We expect that this considerably simpler approach will allow for the relaxation of additional model space assumptions in the future. In particular, we plan to use this framework to learn the system timescale causal structure from subsampled data when latent time series confound our observations.

Acknowledgements

We thank the anonymous reviews for comments that improved this paper. AH was supported by Academy of Finland Centre of Excellence in Computational Inference Research COIN (grant 251170) and Academy of Finland grant 295673. SP was supported by NSF IIS-1318759 & NIH R01EB005846. MJ was supported by COIN (grant 251170) and Academy of Finland grants 276412, 284591; and Research Funds of the University of Helsinki. FE was supported by NSF 1564330. DD was supported by NSF IIS-1318815 & NIH U54HG008540 (from the National Human Genome Research Institute through funds provided by the trans-NIH Big Data to Knowledge (BD2K) initiative). The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

References

- [1] C. Ansótegui, J. Gabàs, Solving (weighted) partial MaxSAT with ILP, in: C.P. Gomes, M. Sellmann (Eds.), *Integration of AI and OR Techniques in Constraint Programming for Combinatorial Optimization Problems*, 10th International Conference, in: Lect. Notes Comput. Sci., vol. 7874, Springer, 2013, pp. 403–409.
- [2] G. Audemard, L. Simon, Predicting learnt clauses quality in modern SAT solvers, in: *Proceedings of the 21st International Joint Conference on Artificial Intelligence*, 2009, pp. 399–404.
- [3] A. Biere, Splatz, Lingeling, Plingeling, Treengeling, YalSAT entering the SAT competition 2016, in: T. Balyo, M. Heule, M. Jarvisalo (Eds.), *Proc. of SAT Competition 2016—Solver and Benchmark Descriptions*, Department of Computer Science Series of Publications B, vol. B-2016-1, University of Helsinki, 2016, pp. 44–45.
- [4] A. Biere, M. Heule, H. van Maaren, T. Walsh (Eds.), *Handbook of Satisfiability*, FAIA, vol. 185, IOS Press, 2009.

¹¹ This complication is independent of undersampling, and arises even if $u = 1$.

- [5] N. Bjørner, N. Narodytska, Maximum satisfiability using cores and correction sets, in: Q. Yang, M. Wooldridge (Eds.), *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence*, AAAI Press, 2015, pp. 246–252.
- [6] D. Danks, S. Plis, Learning causal structure from undersampled time series, in: *Proceedings of the NIPS 2013 Workshop on Causality*, 2013.
- [7] D. Dash, M. Druzdzel, Caveats for causal reasoning with equilibrium models, in: *Symbolic and Quantitative Approaches to Reasoning with Uncertainty*, 6th European Conference, Springer, in: *Lect. Notes Comput. Sci.*, vol. 2143, 2001, pp. 192–203.
- [8] J. Davies, F. Bacchus, Exploiting the power of MIP solvers in MAXSAT, in: M. Järvisalo, A.V. Gelder (Eds.), *16th International Conference Theory and Applications of Satisfiability Testing*, SAT 2013, in: *Lect. Notes Comput. Sci.*, vol. 7962, Springer, 2013, pp. 166–181.
- [9] D. Entner, P. Hoyer, On causal discovery from time series data using FCI, in: *Proceedings of the 5th European Workshop on Probabilistic Graphical Models*, 2010, pp. 121–128.
- [10] M. Gebser, B. Kaufmann, R. Kaminski, M. Ostrowski, T. Schaub, M. Schneider, Potassco: the Potsdam answer set solving collection, *AI Commun.* 24 (2011) 107–124.
- [11] M. Gong, K. Zhang, B. Schoelkopf, D. Tao, P. Geiger, Discovering temporal causal relations from subsampled data, in: *Proceedings of the 32nd International Conference on Machine Learning*, JMLR.org, *J. Mach. Learn. Res. Workshop Conf. Proc.* 37 (2015) 1898–1906.
- [12] C. Granger, Investigating causal relations by econometric models and cross-spectral methods, *Econometrica* 37 (3) (1969) 424–438.
- [13] C. Granger, Testing for causality: a personal viewpoint, *J. Econ. Dyn. Control* 2 (1980) 329–352.
- [14] C. Granger, Some recent development in a concept of causality, *J. Econom.* 39 (1) (1988) 199–211.
- [15] J. Hamilton, *Time Series Analysis*, vol. 2, Princeton University Press, 1994.
- [16] A. Hyttinen, F. Eberhardt, M. Järvisalo, Constraint-based causal discovery: conflict resolution with answer set programming, in: N.L. Zhang, J. Tian (Eds.), *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence*, AUAI Press, 2014, pp. 340–349.
- [17] A. Hyttinen, S. Plis, M. Järvisalo, F. Eberhardt, D. Danks, Causal discovery from subsampled time series data by constraint optimization, in: A. Antonucci, G. Corani, C.P. de Campos (Eds.), *Probabilistic Graphical Models – Eighth International Conference*, JMLR.org, *J. Mach. Learn. Res. Workshop Conf. Proc.* 52 (2016) 216–227.
- [18] A. Hyvärinen, K. Zhang, S. Shimizu, P. Hoyer, Estimation of a structural vector autoregression model using non-Gaussianity, *J. Mach. Learn. Res.* 11 (2010) 1709–1731.
- [19] Y. Iwasaki, H. Simon, Causality and model abstraction, *Artif. Intell.* 67 (1) (1994) 143–194.
- [20] M. Koshimura, T. Zhang, H. Fujita, R. Hasegawa, QMaxSAT: a partial max-sat solver, *J. Satisf. Boolean Model. Comput.* 8 (1/2) (2012) 95–100.
- [21] M. Kutz, The complexity of Boolean matrix root computation, *Theor. Comput. Sci.* 325 (3) (2004) 373–390.
- [22] H. Lütkepohl, *New Introduction to Multiple Time Series Analysis*, Springer Science & Business Media, 2005.
- [23] S. Magliacane, T. Claassen, J.M. Mooij, Ancestral causal inference, in: D.D. Lee, M. Sugiyama, U.V. Luxburg, I. Guyon, R. Garnett (Eds.), *Adv. Neural Inf. Process. Syst.*, vol. 29, Curran Associates, Inc., 2016, pp. 4466–4474.
- [24] D. Margaritis, F. Bromberg, Efficient Markov network discovery using particle filters, *Comput. Intell.* 25 (4) (2009) 367–394.
- [25] R. Martins, V.M. Manquinho, I. Lynce, Open-WBO: a modular MaxSAT solver, in: C. Sinz, U. Egly (Eds.), *Theory and Applications of Satisfiability Testing*, SAT 2014, 17th International Conference, in: *Lect. Notes Comput. Sci.*, vol. 8561, Springer, 2014, pp. 438–445.
- [26] A. Morgado, A. Ignatiev, J. Marques-Silva, MSCG: robust core-guided MaxSAT solving, *J. Satisf. Boolean Model. Comput.* 9 (2015) 129–134.
- [27] N. Narodytska, F. Bacchus, Maximum satisfiability using core-guided maxsat resolution, in: C.E. Brodley, P. Stone (Eds.), *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence*, AAAI Press, 2014, pp. 2717–2723.
- [28] I. Niemelä, Logic programs with stable model semantics as a constraint programming paradigm, *Ann. Math. Artif. Intell.* 25 (3–4) (1999) 241–273.
- [29] J. Peters, D. Janzing, B. Schölkopf, Causal inference on time series using restricted structural equation models, in: C.J.C. Burges, L. Bottou, M. Welling, Z. Ghahramani, K.Q. Weinberger (Eds.), *Adv. Neural Inf. Process. Syst.*, vol. 26, Curran Associates, Inc., 2013, pp. 154–162.
- [30] S. Plis, D. Danks, C. Freeman, V. Calhoun, Rate-agnostic (causal) structure learning, in: *Adv. Neural Inf. Process. Syst.*, vol. 28, Curran Associates, Inc., 2015, pp. 3285–3293.
- [31] S. Plis, D. Danks, J. Yang, Mesochronal structure learning, in: *Proceedings of the 31st Conference on Uncertainty in Artificial Intelligence*, AUAI Press, 2015, pp. 702–711.
- [32] P. Saikko, J. Berg, M. Järvisalo, LMHS: a SAT-IP hybrid MaxSAT solver, in: N. Creignou, D.L. Berre (Eds.), *Theory and Applications of Satisfiability Testing*, SAT 2016, 19th International Conference, in: *Lect. Notes Comput. Sci.*, vol. 9710, Springer, 2016, pp. 539–546.
- [33] P. Simons, I. Niemelä, T. Soinen, Extending and implementing the stable model semantics, *Artif. Intell.* 138 (1–2) (2002) 181–234.
- [34] D. Sonntag, M. Järvisalo, J. Peña, A. Hyttinen, Learning optimal chain graphs with answer set programming, in: *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence*, AUAI Press, 2015, pp. 822–831.
- [35] P. Spirtes, C. Glymour, R. Scheines, *Causation, Prediction, and Search*, Springer, 1993.
- [36] A. Tank, E. Fox, A. Shojaie, Identifiability of non-Gaussian structural VAR models for subsampled and mixed frequency time series, in: *The 2016 ACM SIGKDD Workshop on Causal Discovery*, 2016.
- [37] W. Wei, *Time Series Analysis*, Addison-Wesley, 1994.