

A Constraint Optimization Approach to Causal Discovery from Subsampled Time Series Data

Antti Hyttinen¹

HIIT, Department of Computer Science, University of Helsinki

Sergey Plis

Mind Research Network and University of New Mexico

Matti Järvisalo

HIIT, Department of Computer Science, University of Helsinki

Frederick Eberhardt

Humanities and Social Sciences, California Institute of Technology

David Danks

Department of Philosophy, Carnegie Mellon University

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

Email addresses: `antti.hyttinen@helsinki.fi` (Antti Hyttinen), `s.m.plis@gmail.com` (Sergey Plis), `matti.jarvisalo@helsinki.fi` (Matti Järvisalo), `fde@caltech.edu` (Frederick Eberhardt), `ddanks@cmu.edu` (David Danks)

¹Corresponding author

of non-parametric estimation of system timescale causal structures from subsampled time series data.

Keywords: causality, causal discovery, graphical models, time series, constraint satisfaction, constraint optimization.

1. Introduction

Time-series data has long constituted the basis for causal modeling in many fields of science (Granger, 1969; Hamilton, 1994; Lütkepohl, 2005). 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 directionality, but time series data that undersamples the generating process can be especially misleading about the true direct causal connections (Dash and Druzdzel, 2001; Iwasaki and Simon, 1994).

For example, Figure 1a shows the causal structure of a process unrolled over discrete time steps, and Figure 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 Figure 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 Figure 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 (Granger, 1969, 1980)) 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 (Lütkepohl, 2005; Hamilton, 1994; Hyvärinen et al., 2010)), though only very specific types of interactions can be captured with these latter models. In contrast, we follow Plis et al. (2015a,b) and Gong et al. (2015), and explore the possibility of identifying (features of) the causal process at the true timescale from data that subsample this process.

Plis et al. (2015a,b) 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

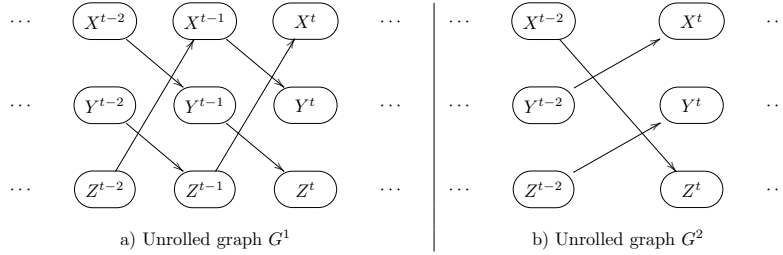


Figure 1: (a) The structure of the causal system-scale time series. (b) The structure of the corresponding measurement scale time series if only every second sample is observed i.e. nodes at time slice $t - 1$ are marginalized. If subsampling is ignored and (b) is thought to depict the true causal structure,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

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. (2015) 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 (Biere et al., 2009; Gebser et al., 2011), and demonstrate that it is orders of magnitudes faster than the current state-of-the-art method by Plis et al. (2015b). 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. (2015), 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) (Hytinen et al., 2016). 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 (Entner and Hoyer, 2010). 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. (2015b), 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 $X^t \rightarrow Y^t$ (Granger, 1988). 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 Spirtes et al. (1993)) 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 Figure 2a for an example). Since the causal process is time-invariant, the edges repeat through t . In accordance with Plis et al. (2015b), 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. Figure 2b shows the rolled graph representation \mathcal{G}^1 of G^1 in Figure 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 (Gong et al., 2015), 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. Figure 2c shows the measurement timescale structure G^2 corresponding to subsampling rate $u = 2$ for the system timescale causal structure in Figure 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

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

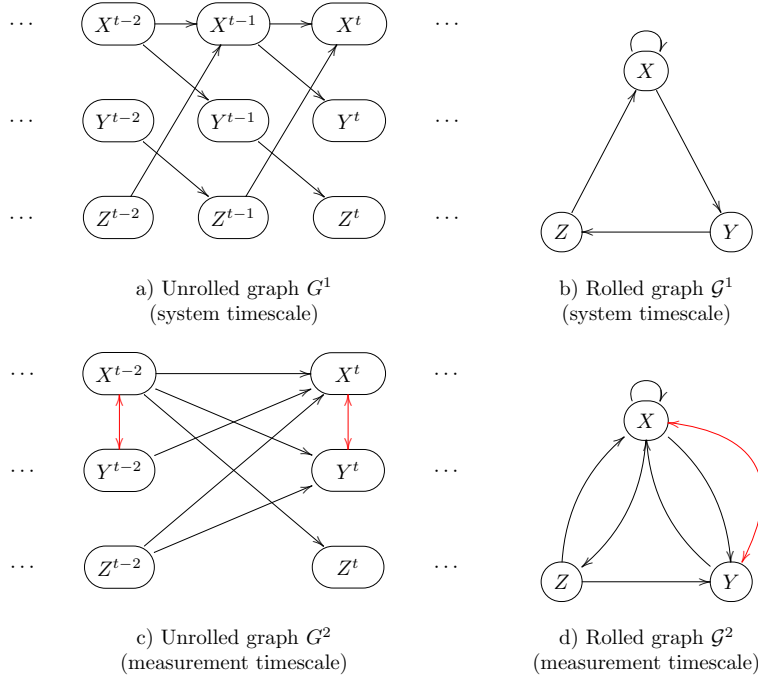


Figure 2: Graph (a) shows the unrolled system timescale structure, where edges repeat through time steps. Graph (b) shows the rolled representation of the same structural information. Graph (c) shows the measurement timescale structure for subsampling rate $u = 2$, i.e. nodes at time slice $t - 1$ in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

where \rightsquigarrow denotes a directed path.

G^u must also represent “direct” connections between variables in the same time step (Wei, 1994). The bi-directed arrow $X^t \leftrightarrow Y^t$ in Figure 2c is an example: X^{t-1} is an unobserved (in the data) common cause of X^t and Y^t in G^1 (Figure 2a). Formally, the system timescale structure G^1 induces bi-directed edges in the measurement timescale G^u as follows:

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

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:

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

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

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:

$$\begin{aligned} X \rightarrow Y \in \mathcal{G}^{u \cdot k} &\Leftrightarrow X \overset{u}{\rightsquigarrow} Y \in \mathcal{G}^k, \\ X \leftrightarrow Y \in \mathcal{G}^{u \cdot k} &\Leftrightarrow \exists Z, l < u : (X \overset{l}{\rightsquigarrow} Z \overset{l}{\rightsquigarrow} Y) \in \mathcal{G}^k \quad \vee \\ &\quad \exists Z, W, l < u : (X \overset{l}{\rightsquigarrow} Z \leftrightarrow W \overset{l}{\rightsquigarrow} Y) \in \mathcal{G}^k, \quad \text{where } X \neq Y. \end{aligned}$$

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:

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

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

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 (2013) 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 $\langle \mathcal{G}^1, u \rangle$ 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*

161 structures \mathcal{G}^1 (at the system timescale) that are (optimally) consistent with the
 162 data.

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

166 3. Finding Consistent System Timescale Structures

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

171 3.1. On Computational Complexity

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

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

177 *Proof.* Membership in NP follows from a guess and check: guess a candidate
 178 \mathcal{G}^1 , and deterministically check whether the length- u paths of \mathcal{G}^1 correspond to
 179 the edges of \mathcal{H} (Plis et al., 2015b). For NP-hardness, for any fixed $u \geq 2$, there
 180 is a straightforward reduction from the NP-complete problem of determining
 181 whether a Boolean B matrix² has a u th root (Kutz, 2004): for a given $n \times n$
 182 Boolean matrix B , interpret B as the directed edge relation of \mathcal{H} , i.e., \mathcal{H} has
 183 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
 184 consistent with the obtained \mathcal{H} iff $B = A^u$ for some binary matrix A (i.e., a u th
 185 root of B). \square

186 If u is unknown, then membership in NP can be established in the same
 187 way by guessing both a candidate \mathcal{G}^1 and a value for u . Theorem 1 ignores
 188 the possible bi-directed edges in \mathcal{H} (whose presence/absence is also harder to
 189 determine reliably from practical numbers of samples; see Section 5). Knowledge
 190 of the presences and absences of such edges in \mathcal{H} can restrict the set of candidate
 191 \mathcal{G}^1 s. For example, in the special case where \mathcal{H} is known to not contain *any*
 192 bi-directed edges, the possible \mathcal{G}^1 s have a fairly simple structure: in any \mathcal{G}^1
 193 that is consistent with \mathcal{H} , every node has at most one successor.³ Whether this
 194 knowledge can be used to prove a more fine-grained complexity result for special
 195 cases is an open question.

²Multiplication of two values in $\{0, 1\}$ is defined as the logical-or, or equivalently, the maximum operator.

³To 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.

3.2. 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. (2015b); 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. (2015b), 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 formalism⁴ (Niemelä, 1999; Simons et al., 2002; Gebser et al., 2011). 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 (Gebser et al., 2011) 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
```

⁴Note the comparison to other solvers using the propositional SAT formalism in Section 7.

235 Here the *cardinality constraint* $1 \{ u(U) : \text{urange}(U) \} 1$ states that the pred-
 236 icate u is true for exactly one value U chosen from those for which $\text{urange}(U)$ is
 237 true.

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

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

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

247 The implied measurement timescale structure \mathcal{G}^u for a candidate solution \mathcal{G}^1
 248 is represented using the predicates $\text{edgeu}/2$ and $\text{confu}/2$, which are derived in the
 249 following way. First, we declare the mapping from a given \mathcal{G}^1 to the correspond-
 250 ing \mathcal{G}^u by declaring the exact length- L paths in a non-deterministically chosen
 251 candidate solution \mathcal{G}^1 . For this, we declare rules that compute the length- L
 252 paths inductively for all $L \leq U$, using the predicate $\text{path}(X,Y,L)$ to represent
 253 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).
```

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

258 Second, to obtain \mathcal{G}^u , we encode Equations 1 and 2 with the following rules
 259 that form predicates edgeu and confu describing the edges \mathcal{G}^1 induces on the
 260 measurement timescale structure \mathcal{G}^u . The first rule derives induced directed
 261 edges in \mathcal{G}^u from the length- U paths, and the second the bidirected edges based
 262 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).
```

263 Finally, we declare constraints that require that the \mathcal{G}^u represented by the
 264 edgeu and confu predicates is consistent with the input \mathcal{H} . This is achieved with
 265 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. (2015b) are correct and complete, so we focus on differences in efficiency, using the implementation of MSL by the original authors. Our approach allows for searching simultaneously over a range of values of u , but Plis et al. (2015b) 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

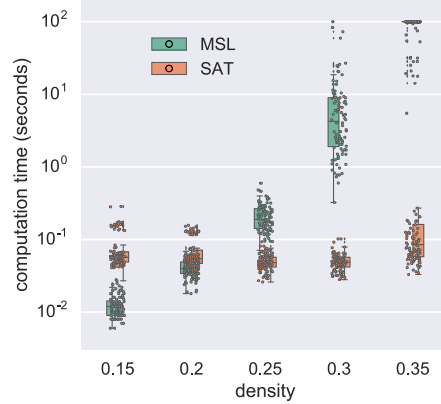


Figure 3: Running times for 10-node rolled graphs as a function of graph density for the state of the art (MSL) and our method (SAT). We used 100 graphs per density and a timeout of 100 seconds; both methods enumerate up to 1000 solutions.

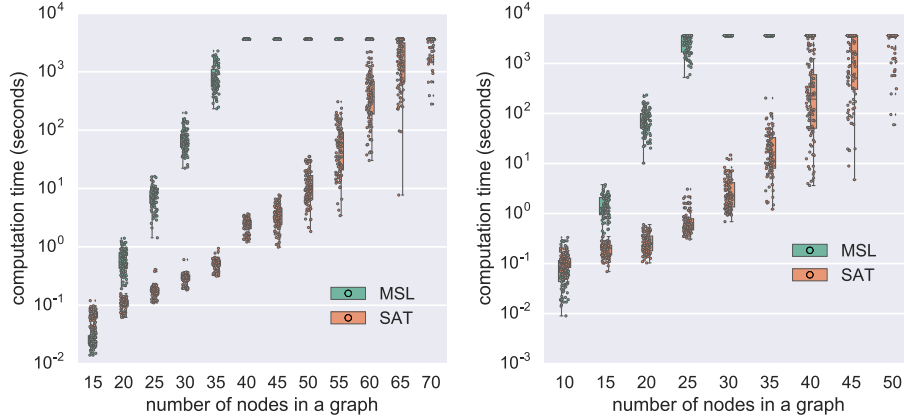


Figure 4: Running times as function of the number of nodes for the state of the art (MSL) and our method (SAT). Left: 10%-dense graphs. Right: 15%-dense graphs. In both plots we use 100 graphs per size and a timeout of 1 hour; both methods enumerate up to 1000 solutions.

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 Figure 3. Figure 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

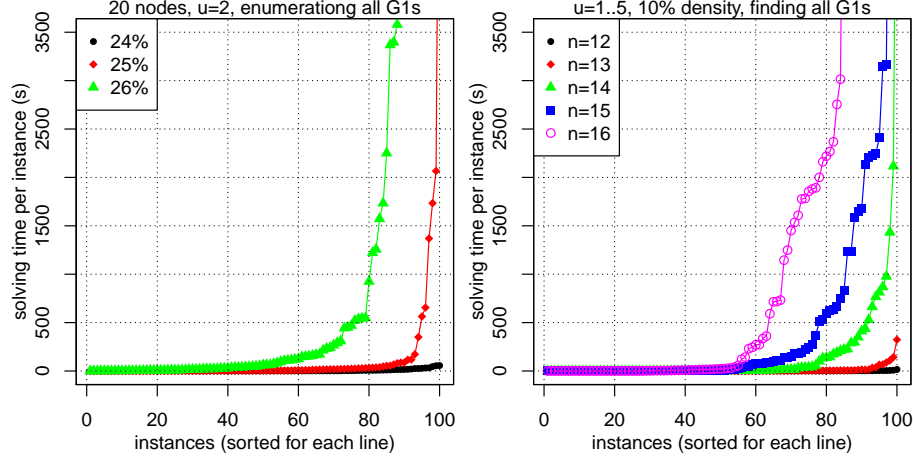


Figure 5: Left: Influence of input graph density on running times of our approach when the subsampling rate $u = 2$ is given as input and all solutions are enumerated. Right: Scalability of our approach when u is left to be determined by the method from interval $1, \dots, 5$. All solutions over the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

313 nodes with a one hour time limit; even for 70 nodes, 25 graphs finished in one
 314 hour. In contrast, MSL reaches only 35 nodes; our approach uses only a few sec-
 315 onds for those graphs. The scalability of our algorithm allows for investigating
 316 the influence of edge density for larger graphs. Figure 5 (left) plots the running
 317 times of our approach (when enumerating *all* solutions) for $u = 2$ ($u = 2$ was
 318 given as input) on 20-node input graphs of varying densities. Note that here
 319 the instances are sorted by the running time for each individual density (curve).
 320 With a time limit of 1000 seconds we can solve 80% of the instances with 26%
 321 density, almost all of the instances with 25% density and all of the instances
 322 with 24% density. Thus, the running time is increased for denser graphs: in
 323 addition to more constraints, there are also more members in the equivalence
 324 classes. Finally, Figure 5 (right) shows the scalability of our approach in the
 325 more challenging task of enumerating *all* solutions over the *range* $u = 1, \dots, 5$
 326 simultaneously. This also demonstrates the generality of our approach: it is not
 327 restricted to solving for individual values of u separately.

328 4. Learning System Timescale Structures from Data

329 Due to statistical errors in estimating \mathcal{H} and the sparse distribution of im-
 330 plied \mathcal{G}^u in the space of possible undersampled graphs, the estimated \mathcal{H} will
 331 often have *no* \mathcal{G}^1 s with $\mathcal{G}^u = \mathcal{H}$. Given such an \mathcal{H} , neither the MSL algorithm
 332 nor our approach in the previous section can output a solution, and they simply

333 conclude that no solution \mathcal{G}^1 exists for the input \mathcal{H} .⁵ In terms of our constraint
 334 declarations, this is witnessed by conflicts among the constraints and the under-
 335 lying model space for any possible solution candidate. Given the inevitability
 336 of statistical errors, we should not simply conclude that no consistent \mathcal{G}^1 ex-
 337 ists for such an \mathcal{H} . Rather, we should aim to learn \mathcal{G}^1 s that, in light of the
 338 underlying conflicts, are “optimally close” (in some well-defined sense of “opti-
 339 mal”) to being consistent with \mathcal{H} . We now turn to this more general problem
 340 setting, and propose what (to the best of our knowledge) is the first approach
 341 to learning, by employing constraint optimization, from undersampled data un-
 342 der conflicts. In fact, we can use the ASP formulation already discussed—with
 343 minor modifications—to address this problem.

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

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

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

358 4.1. Learning by Constraint Optimization

359 To model the objective function for handling conflicts, only simple modifi-
 360 cations are needed to our ASP encoding: instead of declaring *hard* constraints
 361 that require that the paths induced by \mathcal{G}^1 *exactly* correspond to the edges in
 362 \mathcal{H} , we *soften* these constraints by declaring that the violation of each individual
 363 constraint incurs the associated weight as penalty. In the ASP language, this
 364 can be expressed by augmenting the input predicates `edgeh(X,Y)` with weights:
 365 `edgeh(X,Y,W)` (and similarly for `no_edgeh`, `confh` and `no_confh`), and by using
 366 *weighted soft rules* syntactically represented via `:~` instead of `:-`. Here the ad-
 367 ditional argument W represents the weight $w((X \rightarrow Y) \in \mathcal{H})$ given as input.
 368 The following expresses that each conflicting presence of an edge in \mathcal{H} and \mathcal{G}^u is
 369 penalized with the associated weight W . The additional `[W,X,Y,v]` for $v = 1, 2$

⁵For these cases, Plis et al. (2015b) 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.

370 syntactically enforce that a cost of W is incurred in case the corresponding rule
 371 is violated for a specific pair of nodes X, Y . The numbers $v \in \{1, 2\}$ at the
 372 end of the brackets enable the solver to distinguish the cost incurred due to
 373 bidirected and directed edges respectively.

```

:~ edgeh(X,Y,W), not edgeu(X,Y). [W,X,Y,1]
:~ no_edgeh(X,Y,W), edgeu(X,Y). [W,X,Y,1]
:~ confh(X,Y,W), not confu(X,Y). [W,X,Y,2]
:~ no_confh(X,Y,W), confu(X,Y). [W,X,Y,2]

```

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

377 4.2. Weighting Schemes

378 We use two different schemes for weighting the presences and absences of
 379 edges in \mathcal{H} according to their reliability. To determine the presence or absence
 380 of a specific edge $X \rightarrow Y$ in \mathcal{H} , we simply test the corresponding independence
 381 $X^{t-1} \perp\!\!\!\perp Y^t \mid \mathbf{V}^{t-1} \setminus X^{t-1}$. To determine the presence/absence of an edge $X \leftrightarrow Y$
 382 in \mathcal{H} , we test the independence: $X^t \perp\!\!\!\perp Y^t \mid \mathbf{V}^{t-1}$.

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

$$\begin{aligned}
 w(e \in \mathcal{H}) &= 1 \quad \forall e \in \mathcal{H}, \\
 w(e \notin \mathcal{H}) &= 1 \quad \forall e \notin \mathcal{H}.
 \end{aligned}$$

384 Uniform edge weights resemble the search on the Hamming cube of \mathcal{H} that
 385 Plis et al. (2015b) used to address the problem of finding \mathcal{G}^1 s when \mathcal{H} did not
 386 correspond to any \mathcal{G}^u , though our approach is much superior computationally.

387 A more intricate approach is to use pseudo-Bayesian weights following Mar-
 388 garitis and Bromberg (2009); Hyttinen et al. (2014); Sonntag et al. (2015).
 389 They used Bayesian model selection to obtain reliability weights for indepen-
 390 dence tests. Instead of a p -value and a binary decision, these types of tests give
 391 a measurement of reliability for an independence/dependence statement as a
 392 Bayesian probability. We can directly incorporate their approach of using log-
 393 probabilities as the reliability weights for the edges. For details, see Section 4.3
 394 of Hyttinen et al. (2014). Again, we only compute weights for the independence
 395 tests mentioned above in the estimation of \mathcal{H} .

396 5. Simulations

397 We use simulations to explore the accuracy and runtime efficiency of our
 398 approach in various different settings. For the simulations, system timescale
 399 structures \mathcal{G}^1 and the associated data generating models were constructed in
 400 the following way. To guarantee connectedness of the graphs, we first formed
 401 a cycle of all nodes in a random order (following Plis et al. (2015b)). We
 402 then randomly sampled additional directed edges until the required density was

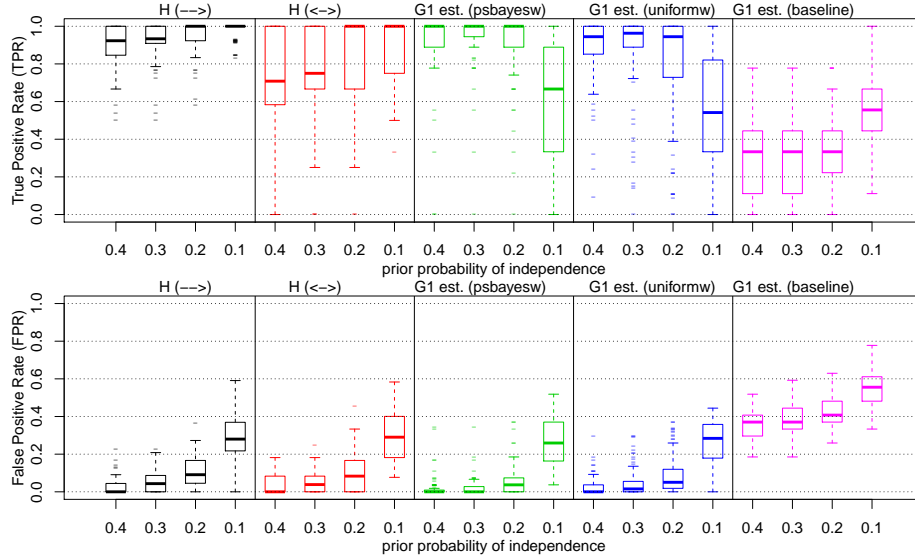


Figure 6: Accuracy of the optimal solutions when subsampling rate $u = 2$ is given as input (200 instances and 250 samples). The x-axis shows the different prior probabilities of independence in the utilized independence test. The two left columns give the accuracy of the estimation of the measurement 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).

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 samples⁶ to get the particular subsampling rate.⁷

5.1. Accuracy

Figure 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

⁶All sample counts refer to the number of samples after subsampling.

⁷Clingo only accepts integer weights; we multiplied weights by 1000 and rounded to the nearest integer.

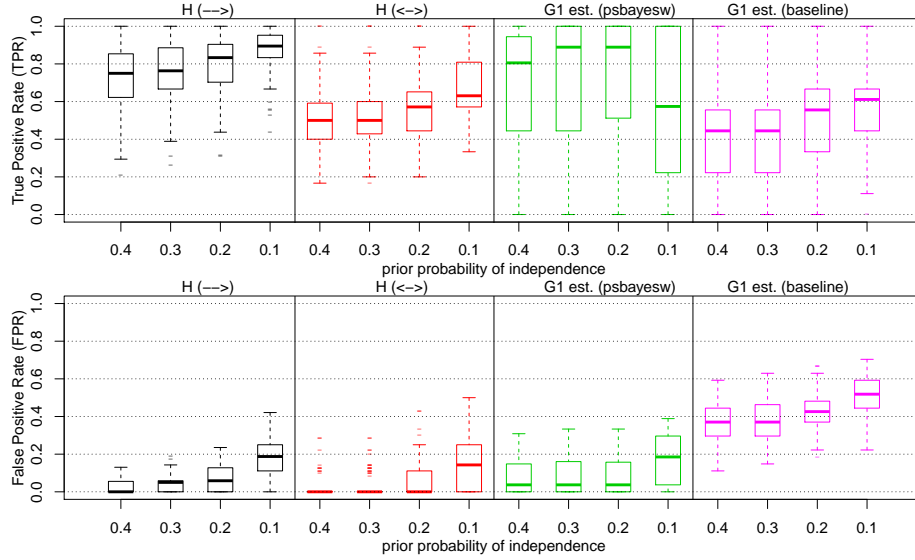


Figure 7: Accuracy of the optimal solutions when subsampling rate $u = 2$ is given as input (200 instances and 500 samples). The x-axis shows the different prior probabilities of independence in the utilized independence test. The two left columns give the accuracy of the estimation of the measurement 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.

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 parameter values, but not perfectly. The presence of directed edges can be estimated more accurately. More importantly, the two rightmost columns in Figure 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.

Figure 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 (Figure 7, black, red). The problem is that measurement timescale edges here correspond to 3-edge paths, whose

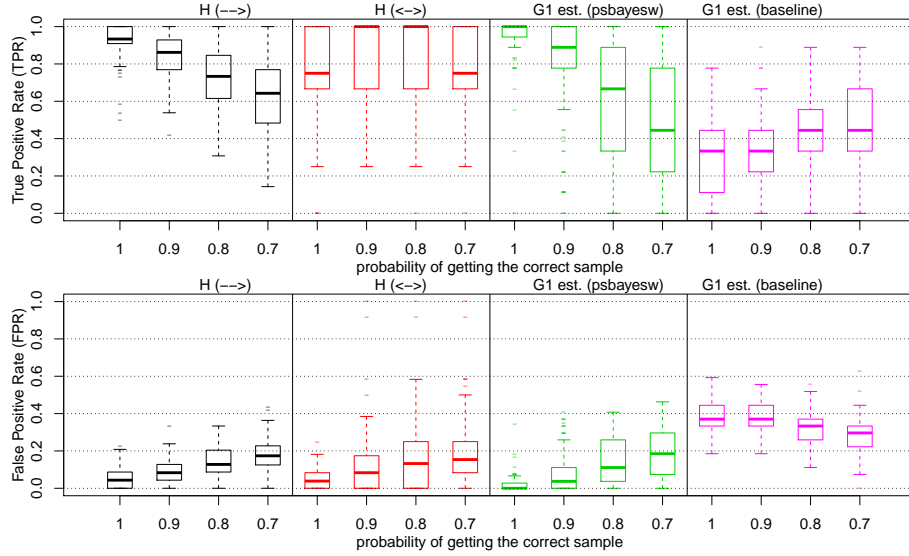


Figure 8: Accuracy of the optimal solutions when subsampling rate $u = 2$ is given as input (200 instances and 250 samples), some samples are obtained at the adjacent timepoints. Due to previous simulations we used the prior probability of 0.3 for all methods. In more detail, the x-axis gives the probability 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.

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 (Figure 7, 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

Figure 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 Figure 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

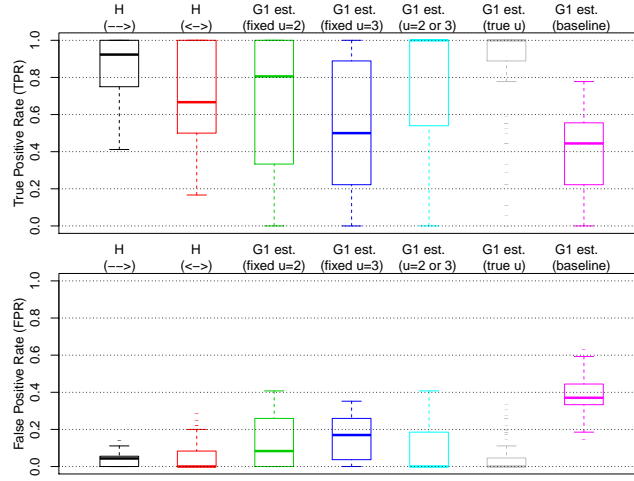


Figure 9: Accuracy when the true u is unknown. Two left boxplots show accuracy of the \mathcal{H} estimate as before. The next three boxplots show the accuracy of our approach (pseudo-Bayesian weights) when, regardless of the true u , u is fixed to 2, or to 3, or left for the procedure decision, respectively. 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.

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.

Figure 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 second from the right (the method was run with $u = 2 \dots 3$). In fact, the accuracy comes close to that of the rightmost boxplots, where the correct u was given as input to the procedure. Thus the procedure

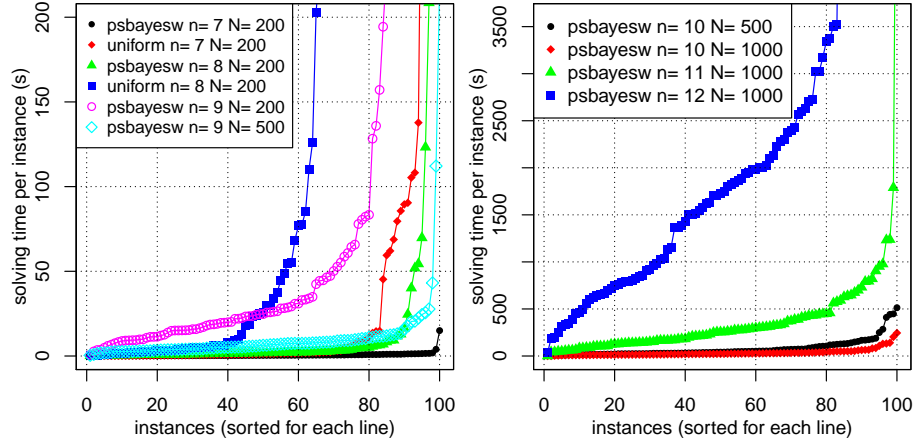


Figure 10: Scalability of our constraint optimization approach (using *Clingo*) for different graph sizes, numbers of samples and weighting schemes. For each setting there are 100 instances that are sorted according to the solving time on each line.

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 Figure 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. Figure 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 line), while the uniform weighting scheme (red line) 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 line) take longer than for $n = 9, N = 500$ (Figure 10 left, magenta vs. cyan lines). 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 (Figure 10 right), though in a considerable amount of time.

6. Case Study: House data of Peters et al. (2013)

In order to demonstrate the applicability to real-world data, we analyzed the house temperature and humidity data of Peters et al. (2013). 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

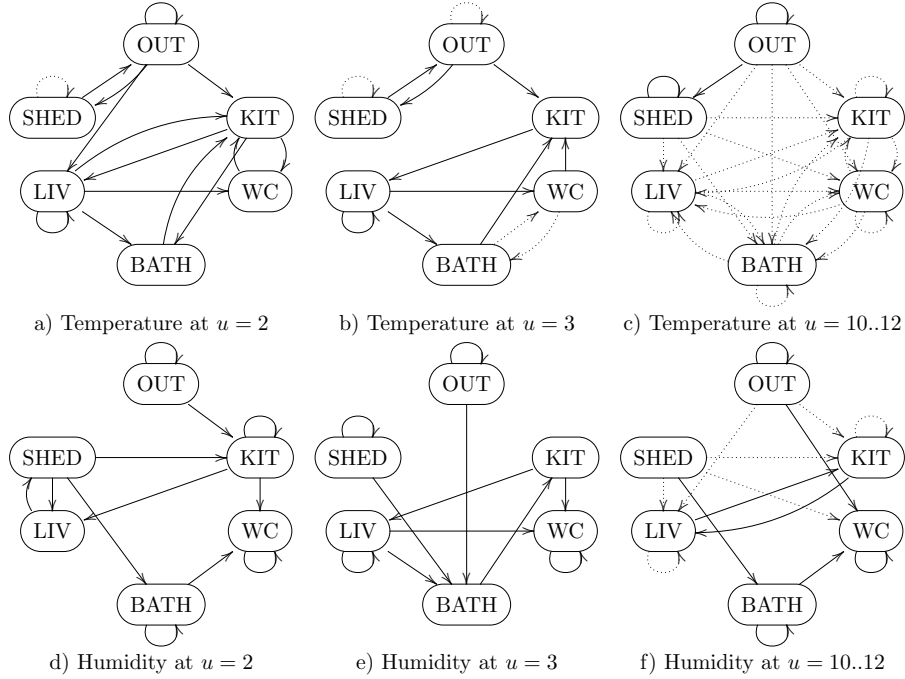


Figure 11: Results of the House data analysis for different subsampling rates (u) and measurement type. Edges with full lines are found to be present, absent edges are found to be absent, edges with dotted lines may be present or absent.

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. (2015). 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.

We analyzed the temperature and humidity components separately, and examined the differences of sequential measurements,⁸ 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 (Magliacane et al., 2016). For each edge in

⁸This may take out some of the influences of self-loops.

\mathcal{G}^1 , we ran the inference procedure first enforcing its presence and then enforcing its absence.⁹ 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 $u = 3$ (Figure 11a,b) differ substantially from one another, as do the two humidity graphs (Figure 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. Figure 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. Figure 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. (2013). Similarly as Peters et al. (2013), 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. (2013). 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.

⁹This 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 $X \rightarrow Y$.

559 7. Solver Performance Comparison

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

570 7.1. Direct Propositional SAT Encoding

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

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

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

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

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

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

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

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

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

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

578 7.2. Solver Comparison: Finding Consistent System Timescale Structures

579 The results of a runtime performance comparison between **Clingo** and two
 580 state-of-the-art SAT solvers, Glucose (Audemard and Simon, 2009) and Lin-
 581 geling (Biere, 2016), is presented in Figure 12 for $u = 3$ (given as input), edge
 582 density of 10% and the numbers of nodes ranging from 27 (on left) to 30 (on

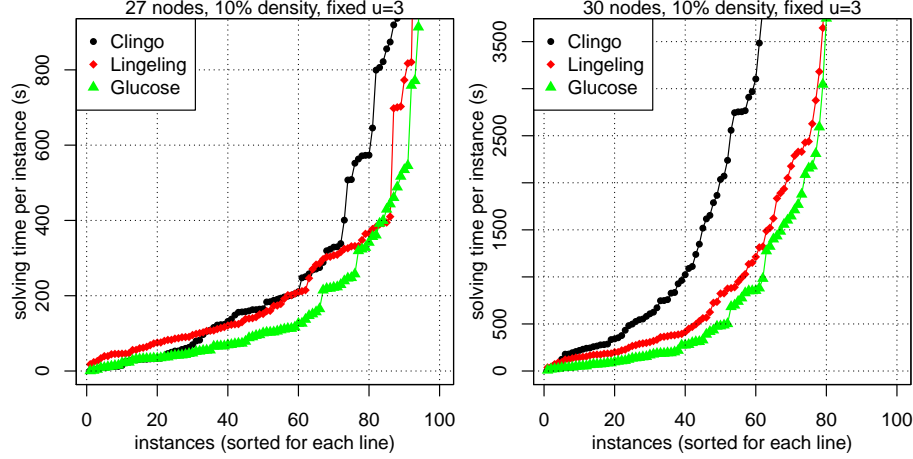


Figure 12: Comparison of running times for different solvers finding a single graph in the equivalence class, when the subsampling rate $u = 3$ is given as input. Left: easier instances with 27 nodes. Right: harder instances with 30 nodes. **Clingo** uses the ASP encoding presented in Section 3.2, all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

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 (Narodytska and Bacchus, 2014), LMHS (Saikko et al., 2016), MSCG (Morgado et al., 2015), Open-WBO (Martins et al., 2014), PrimalDual (Bjørner and Narodytska, 2015), and QMaxSAT (Koshimura et al., 2012), as well as the commercial

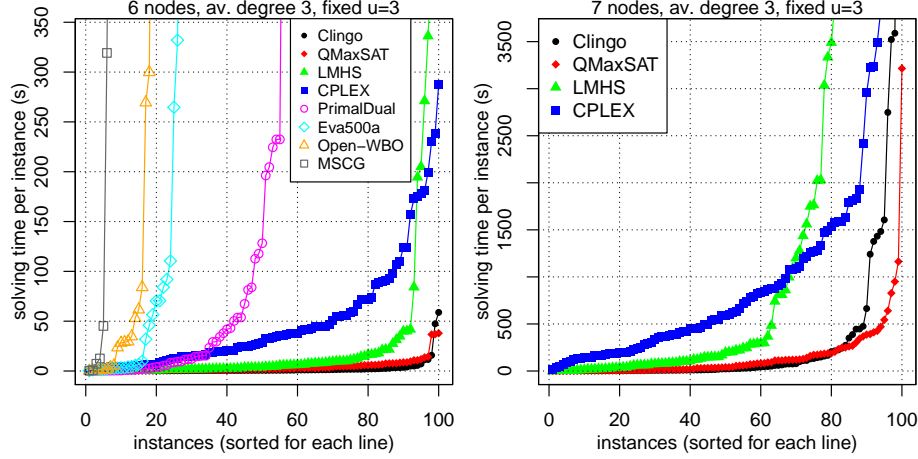


Figure 13: Comparison of running times for different solvers finding the optimal graph, when the subsampling rate $u = 3$ is given as input. Left: easier instances with 6 nodes. Right: harder instances with 7 nodes. `Clingo` uses the ASP encoding presented in Sections 3.2 and 4.1, all others use the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

integer programming (IP) solver CPLEX run on a standard IP translation of MaxSAT (Davies and Bacchus, 2013; Ansótegui and Gabàs, 2013), are shown in Figure 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.¹⁰ For approaches to this case see Tank et al. (2016), 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

¹⁰For 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.

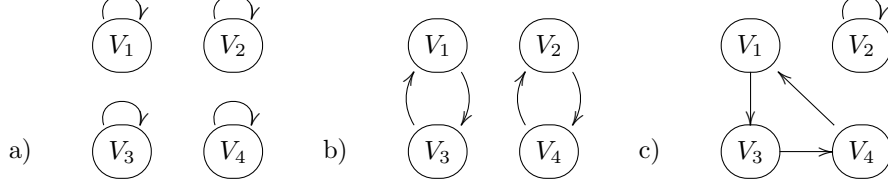


Figure 14: Example graphs for learning from mixed frequency data. Graph (a) shows the true system timescale causal structure. When this is subsampled by $u = 2$ or by $u = 3$, the result is also the structure (a) (this time in measurement timescale). System timescale structure (b) gives measurement timescale 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).

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 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 Figure 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 Figure 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 Figure 14b produces Figure 14a with $u = 2$, and Figure 14c produces Figure 14a with $u = 3$. In fact, *any* system timescale edge can be present or absent given either of the measurement timescale graphs alone.¹¹ However, if this measurement timescale graph is found at *both* $u = 2$ and $u = 3$, then the system timescale structure can be uniquely determined: Figure 14b produces a different measurement timescale graph for $u = 3$ and Figure 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 ($= 2^{4 \cdot 4}$) different 4-variable \mathcal{G}^1 s, and compared the number of equiv-

¹¹The node labels in Figure 14b and c can be permuted.

654 alence classes given input at a single subsampling rate, versus given inputs at
 655 two subsampling rates. A greater number of equivalence classes means a higher
 656 chance that a random graph will be uniquely identifiable, and so the number of
 657 equivalence classes is an approximate (inverse) measure of the extent of under-
 658 determination.

659 For input at a single undersampling rate, for $u = 2$ we have 24265 equiv-
 660 alence classes; 7544 for $u = 3$; and 3964 equivalence classes for $u = 4$. These
 661 results with a single undersampled input graph thus replicate the known result
 662 that underdetermination is a significant problem, and it rapidly worsens as u
 663 increases (Plis et al., 2015a,b).

664 If we instead have measurement timescale graphs for both $u = 2, 3$, then
 665 we have 26720 equivalence classes, which is only slightly more than the number
 666 for $u = 2$ by itself. That is, underdetermination is not substantially reduced
 667 if we additionally measure at $u = 3$ when we already have measurements at
 668 $u = 2$. Similarly, for $u = 3, 4$ we have 7814 equivalence classes; again, there is
 669 a reduction in underdetermination compared to $u = 3$ by itself, but it is quite
 670 small. This analysis assumes that all \mathcal{G}^1 are equally likely, and it is an open
 671 question whether measurements at different undersampling rates would have
 672 more impact for certain classes of \mathcal{G}^1 (e.g., connected graphs).

673 9. Discussion

674 We have assumed that all common causes of measured variables are them-
 675 selves measured, but this assumption is frequently violated in real-world data.
 676 Constraint satisfaction methods have elsewhere been used with success to iden-
 677 tify causal relations in the presence of unobserved common causes or latent
 678 variables (Hyttinen et al., 2014; Magliacane et al., 2016). For time series data,
 679 dropping the assumption of causal sufficiency (in the system timescale) generates
 680 complications. Even if the system timescale process including latent variables
 681 is assumed to be first order Markov, the Markov order of the measurement
 682 timescale (naturally without the latent variables) can be arbitrarily larger.¹²
 683 That is, variables arbitrarily far in the past can (directly, in the measurement
 684 timescale) cause variables at the current timestep. We would thus need to both
 685 enrich the notation for \mathcal{G}^u to encode the time lags of direct causal effects, and
 686 also modify the statistical tests used to estimate these connections.

687 Moreover, there can be more information contained in the pattern of time
 688 lags (i.e., *which* past variables directly cause the present) than is given by the
 689 Markov order of the system. As just one example, suppose $\{X^{t-2}, X^{t-4}, \dots\} \rightarrow$
 690 Y^t . The simplest (in terms of number of latents) structure that explains these
 691 influences (i) has a latent L through which X influences Y (i.e., $X^{t-2} \rightarrow L^{t-1} \rightarrow$
 692 Y^t); and (ii) L is part of a 2-loop with another latent M (i.e., $L^{t-1} \rightarrow M^t$ and
 693 $L^t \leftarrow M^{t-1}$). In contrast, if we have $\{X^{t-2}, X^{t-3}, \dots\} \rightarrow Y^t$, then the simplest
 694 structure has only a single latent L through which X influences Y , but where L

¹²This complication is independent of undersampling, and arises even if $u = 1$.

695 has a self-loop (i.e., $L^{t-1} \rightarrow L^t$). The pattern of time lags for direct causes—in
696 particular, the absence of certain time lags—thus contains information about
697 the number and causal structure of the latent variables. Estimation of this
698 pattern, however, can be quite complex statistically.

699 Subsampled time series data can be also particularly prone to violations
700 of faithfulness. For example, the underlying process unrolled over time may
701 include directed paths over many time steps that do not result in significant
702 statistical dependence in the observed data. In addition, variables observed
703 over subsequent time steps might be almost deterministically related. If $X^{t-1} \approx$
704 X^{t-2} , then conditioning on X^{t-2} may render the statistical dependence through
705 $Y^t \leftarrow X^{t-1} \rightarrow Z^t$ undetectable from any realistic numbers of samples. In the
706 current framework, both of these situations are treated as estimation errors
707 in \mathcal{H} . Further modeling of these complications may help to achieve improved
708 accuracy. Another option could be to develop parametric approaches instead of
709 the non-parametric one presented in this paper.

710 10. Conclusion

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

726 Acknowledgments

727 We thank the anonymous reviews for comments that improved this paper.
728 AH was supported by Academy of Finland Centre of Excellence in Computa-
729 tional Inference Research COIN (grant 251170) and Academy of Finland grant
730 295673. SP was supported by NSF IIS-1318759 & NIH R01EB005846. MJ was
731 supported by COIN (grant 251170) and Academy of Finland grants 276412,
732 284591; and Research Funds of the University of Helsinki. FE was supported
733 by NSF 1564330. DD was supported by NSF IIS-1318815 & NIH U54HG008540
734 (from the National Human Genome Research Institute through funds provided
735 by the trans-NIH Big Data to Knowledge (BD2K) initiative). The content is

736 solely the responsibility of the authors and does not necessarily represent the
737 official views of the National Institutes of Health.

738 References

- 739 Ansótegui C, Gabàs J. Solving (weighted) partial MaxSAT with ILP. In: Gomes
740 CP, Sellmann M, editors. Integration of AI and OR Techniques in Constraint
741 Programming for Combinatorial Optimization Problems, 10th International
742 Conference. Springer; volume 7874 of *Lecture Notes in Computer Science*;
743 2013. p. 403–9.
- 744 Audemard G, Simon L. Predicting learnt clauses quality in modern SAT solvers.
745 In: Proceedings of the 21st International Joint Conference on Artificial Intel-
746 ligence. 2009. p. 399–404.
- 747 Biere A. Splatz, Lingeling, Plingeling, Treengeling, YalSAT entering the SAT
748 Competition 2016. In: Balyo T, Heule M, Jarvisalo M, editors. Proc. of
749 SAT Competition 2016 – Solver and Benchmark Descriptions. University of
750 Helsinki; volume B-2016-1 of *Department of Computer Science Series of Pub-
751 lications B*; 2016. p. 44–5.
- 752 Biere A, Heule M, van Maaren H, Walsh T, editors. Handbook of Satisfiability;
753 volume 185 of *FAIA*. IOS Press; 2009.
- 754 Bjørner N, Narodytska N. Maximum satisfiability using cores and correction
755 sets. In: Yang Q, Wooldridge M, editors. Proceedings of the Twenty-Fourth
756 International Joint Conference on Artificial Intelligence. AAAI Press; 2015.
757 p. 246–52.
- 758 Danks D, Plis S. Learning causal structure from undersampled time series. In:
759 Proceedings of the NIPS 2013 Workshop on Causality; 2013.
- 760 Dash D, Druzdzel M. Caveats for causal reasoning with equilibrium models. In:
761 Symbolic and Quantitative Approaches to Reasoning with Uncertainty, 6th
762 European Conference. Springer; volume 2143 of *LNCS*; 2001. p. 192–203.
- 763 Davies J, Bacchus F. Exploiting the power of MIP solvers in MAXSAT. In:
764 Jarvisalo M, Gelder AV, editors. Theory and Applications of Satisfiability
765 Testing - SAT 2013 - 16th International Conference. Springer; volume 7962
766 of *Lecture Notes in Computer Science*; 2013. p. 166–81.
- 767 Entner D, Hoyer P. On causal discovery from time series data using FCI.
768 In: Proceedings of the 5th European Workshop on Probabilistic Graphical
769 Models. 2010. p. 121–8.
- 770 Gebser M, Kaufmann B, Kaminski R, Ostrowski M, Schaub T, Schneider M.
771 Potassco: The Potsdam answer set solving collection. AI Communications
772 2011;24(2):107–24.

773 Gong M, Zhang K, Schoelkopf B, Tao D, Geiger P. Discovering temporal causal
774 relations from subsampled data. In: Proceedings of the 32nd International
775 Conference on Machine Learning. JMLR.org; volume 37 of *JMLR W&CP*;
776 2015. p. 1898–906.

777 Granger C. Investigating causal relations by econometric models and cross-
778 spectral methods. *Econometrica* 1969;37(3):424–38.

779 Granger C. Testing for causality: a personal viewpoint. *Journal of Economic*
780 *Dynamics and Control* 1980;2:329–52.

781 Granger C. Some recent development in a concept of causality. *Journal of*
782 *Econometrics* 1988;39(1):199–211.

783 Hamilton J. Time series analysis. volume 2. Princeton University Press, 1994.

784 Hyttinen A, Eberhardt F, Järvisalo M. Constraint-based causal discovery: Con-
785 flict resolution with answer set programming. In: Zhang NL, Tian J, editors.
786 Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelli-
787 gence. AUAI Press; 2014. p. 340–9.

788 Hyttinen A, Plis S, Järvisalo M, Eberhardt F, Danks D. Causal discovery from
789 subsampled time series data by constraint optimization. In: Antonucci A,
790 Corani G, de Campos CP, editors. Probabilistic Graphical Models - Eighth
791 International Conference. JMLR.org; volume 52 of *JMLR Workshop and Con-*
792 *ference Proceedings*; 2016. p. 216–27.

793 Hyvärinen A, Zhang K, Shimizu S, Hoyer P. Estimation of a structural vector
794 autoregression model using non-Gaussianity. *Journal of Machine Learning*
795 *Research* 2010;11:1709–31.

796 Iwasaki Y, Simon H. Causality and model abstraction. *Artificial Intelligence*
797 1994;67(1):143–94.

798 Koshimura M, Zhang T, Fujita H, Hasegawa R. Qmaxsat: A partial max-
799 sat solver. *Journal of Satisfiability, Boolean Modeling and Computation*
800 2012;8(1/2):95–100.

801 Kutz M. The complexity of Boolean matrix root computation. *Theoretical*
802 *Computer Science* 2004;325(3):373–90.

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

805 Magliacane S, Claassen T, Mooij JM. Ancestral causal inference. In: Lee DD,
806 Sugiyama M, Luxburg UV, Guyon I, Garnett R, editors. Advances in Neural
807 Information Processing Systems 29. Curran Associates, Inc.; 2016. p. 4466–74.

808 Margaritis D, Bromberg F. Efficient Markov network discovery using particle
809 filters. *Computational Intelligence* 2009;25(4):367–94.

810 Martins R, Manquinho VM, Lynce I. Open-WBO: A modular MaxSAT solver.
811 In: Sinz C, Egly U, editors. Theory and Applications of Satisfiability Testing
812 - SAT 2014 - 17th International Conference. Springer; volume 8561 of *Lecture*
813 *Notes in Computer Science*; 2014. p. 438–45.

814 Morgado A, Ignatiev A, Marques-Silva J. MSCG: Robust core-guided MaxSAT
815 solving. *Journal on Satisfiability, Boolean Modeling and Computation*
816 2015;9:129–34.

817 Narodytska N, Bacchus F. Maximum satisfiability using core-guided maxsat res-
818 olution. In: Brodley CE, Stone P, editors. Proceedings of the Twenty-Eighth
819 AAAI Conference on Artificial Intelligence. AAAI Press; 2014. p. 2717–23.

820 Niemelä I. Logic programs with stable model semantics as a constraint program-
821 ming paradigm. *Annals of Mathematics and Artificial Intelligence* 1999;25(3-
822 4):241–73.

823 Peters J, Janzing D, Schölkopf B. Causal inference on time series using restricted
824 structural equation models. In: Burges CJC, Bottou L, Welling M, Ghahra-
825 mani Z, Weinberger KQ, editors. Advances in Neural Information Processing
826 Systems 26. Curran Associates, Inc.; 2013. p. 154–62.

827 Plis S, Danks D, Freeman C, Calhoun V. Rate-agnostic (causal) structure learn-
828 ing. In: Advances in Neural Information Processing Systems 28. Curran As-
829 sociates, Inc.; 2015a. p. 3285–93.

830 Plis S, Danks D, Yang J. Mesochronal structure learning. In: Proceedings of the
831 31st Conference on Uncertainty in Artificial Intelligence. AUAI Press; 2015b.
832 p. 702–11.

833 Saikko P, Berg J, Järvisalo M. LMHS: A SAT-IP hybrid MaxSAT solver. In:
834 Creignou N, Berre DL, editors. Theory and Applications of Satisfiability Test-
835 ing - SAT 2016 - 19th International Conference. Springer; volume 9710 of
836 *Lecture Notes in Computer Science*; 2016. p. 539–46.

837 Simons P, Niemelä I, Soinen T. Extending and implementing the stable model
838 semantics. *Artificial Intelligence* 2002;138(1-2):181–234.

839 Sonntag D, Järvisalo M, Peña J, Hyttinen A. Learning optimal chain graphs
840 with answer set programming. In: Proceedings of the Thirty-First Conference
841 on Uncertainty in Artificial Intelligence. AUAI Press; 2015. p. 822–31.

842 Spirtes P, Glymour C, Scheines R. Causation, prediction, and search. Springer,
843 1993.

844 Tank A, Fox E, Shojaie A. Identifiability of non-Gaussian structural VAR
845 models for subsampled and mixed frequency time series. In: The 2016 ACM
846 SIGKDD Workshop on Causal Discovery; 2016.

847 Wei W. Time series analysis. Addison-Wesley, 1994.