
Data Generating Process to Evaluate Causal Discovery Techniques for Time Series Data

Andrew R. Lawrence^{*†} Marcus Kaiser[†] Rui Sampaio Maksim Sipos
causalens, London, UK
{andrew,marcus,ru,max}@causalens.com

Abstract

Going beyond correlations, the understanding and identification of causal relationships in observational time series, an important subfield of Causal Discovery, poses a major challenge. The lack of access to a well-defined ground truth for real-world data creates the need to rely on synthetic data for the evaluation of these methods. Existing benchmarks are limited in their scope, as they either are restricted to a “static” selection of data sets, or do not allow for a granular assessment of the methods’ performance when commonly made assumptions are violated. We propose a flexible and simple to use framework for generating time series data, which is aimed at developing, evaluating, and benchmarking time series causal discovery methods. In particular, the framework can be used to fine tune novel methods on vast amounts of data, without “overfitting” them to a benchmark, but rather so they perform well in real-world use cases. Using our framework, we evaluate prominent time series causal discovery methods and demonstrate a notable degradation in performance when their assumptions are invalidated and their sensitivity to choice of hyperparameters. Finally, we propose future research directions and how our framework can support both researchers and practitioners.

1 Introduction

The aim of *Causal Discovery* is to identify causal relationships from purely observational data. Special interest lies in identifying causal effects for time series data, where individual observations arrive ordered in time. Compared to the case of independent and identically distributed (IID) data, a robust analysis of time series data requires one to address additional difficulties and guard against pitfalls. These difficulties include non-stationarity, which can materialize in shifts in distribution (e.g., a shift in the mean or a higher moment, potentially from interventions). Moreover, real-world time series tend to show various levels of autocorrelation. These can both carry valuable long-term information, but also invalidate typical assumptions for statistical procedures, such as the independence of samples (cf. the Gauss-Markov Theorem for linear regression [7, Chapter 3.3.2]).

One major benefit is that the order of time can help distinguish cause from effect. As the future cannot affect the past, the causal driver can be identified as the variable that occurred first. This is a valid assumption for high-resolution data; however, for lower resolutions, one must consider *contemporaneous* or *instantaneous effects*, where one variable has a causal effect on another variable observed at the same point in time.

Additional complications arise when it comes to evaluating and comparing the performance of individual methods. Frequently, new techniques are evaluated against their own synthetic benchmarks,

^{*}Corresponding author.

[†]These authors contributed equally to this work.

rather than following one “gold standard” such as those which have been established in other domains within machine learning, e.g. MNIST [15, 16] and CIFAR-10/100 [14] for image classification. The lack of a general benchmark with known ground truth makes it difficult to compare individual methods, especially when there is not a publicly available implementation of the new method.

For real-world problems, there is often no known ground truth causal structure and it is impossible to observe all the variables to ensure causal sufficiency [35]. In many cases it is unclear how a causal structure can be defined. For example, consider a system with many highly correlated time series; in such a setup, it is often not straightforward to identify whether an observed effect stems from a single time series, a subset, or all of them. Moreover, real-world data often violate assumptions made for causal discovery methods (such as IID data, or linear relationships between variables). Therefore, it is important to test how individual methods perform when these assumptions are not satisfied.

We propose a flexible, yet easy to use synthetic data generation process based on structural causal models [23], which can be used to benchmark causal discovery methods under a wide-range of conditions. We show the performance of prominent methods when key assumptions are invalidated and demonstrate the sensitivity to the choice of hyperparameter values.

2 Background

2.1 Brief overview of causal discovery methods

We here give a high level overview of Causal Graph Discovery methods, which aim to identify a *Directed Acyclic Graph (DAG)* from purely observational data. A DAG consists of nodes connected by directed edges/links from parent nodes to child nodes. Nodes in the DAG represent the variables in the data and edges indicate direct causes, i.e., a variable is said to be a direct cause of another if the former is a parent of the latter in the DAG. Classical methods for Causal Discovery for IID data either depend on *Conditional Independence Tests* or are based on *Functional Causal Models*. For a recent review of the topic, we refer the reader to [9].

Methods that depend on *Conditional Independence Tests* can be seen as a special case of discovery methods for *Bayesian Networks* [34]. These methods use a series of conditional independence tests to construct a graph that satisfies the *Causal Markov Condition*, cf. [35, Chapter 3.4.1], [34]. Typically, the convergence to a true DAG cannot be guaranteed and the resulting graph is only partially directed (a *Completed Partially DAG (CPDAG)*), cf. [4, 5] for more details). Two subclasses of these algorithms are *Constraint Based Methods* (e.g., PC [34] and (R)FCI [5]) and *Score Based Methods*, which optimize a score that results in a graph that is (as close as possible to) a DAG (e.g., Greedy Equivalence Search (GES), see [4, 20]). Both constraint and score based methods can be combined to obtain *Hybrid Methods*, such as GFCI [9], which combines GES with FCI.

Functional Causal Models prescribe a specific functional form to the relation between variables (see § 2.2). A well-known example is the Linear Non-Gaussian Additive Model (LiNGAM) [31, 32]. A more recent example within this class is NoTEARS [40], which encodes a DAG-constraint as part of a differentiable loss function. There are non-linear extensions of NoTEARS [39, 41] and similar ideas with optimization of a loss for learning DAGs have been used in [21] and [37]. This class of methods returns a functional representation, from which a DAG can be obtained. Note that the latter is typically not assumed to be “causal” as it may not satisfy the Causal Markov Condition.

Time series causal discovery For time series causal discovery there are two notions of a causal graph —the *Full Time Graph (FG)* and the *Summary Graph (SG)*; both defined in [25, Chapter 10.1]. The FG is a DAG whose nodes represent the variables at each point in time, with the convention that future values cannot be parents of present or past values. The SG is a “collapsed” version of the FG, where each node represents a whole time series. There exists an edge $X_i \rightarrow X_j$ in the SG if and only if there exists $t \leq t'$ s.t. $X_i(t) \rightarrow X_j(t')$ in the FG.

The classical approach to causal discovery for time series is *Granger Causality* [10]. Intuitively, for two time series X and Y in a universe of observed time series U , we say that “ X Granger causes Y ” if excluding historical values of X from the universe U decreases the forecasting performance of Y . Non-linear versions of Granger Causality [18] have been proposed and Granger Causality is closely linked to (the non-linear) *Transfer Entropy*, cf. [30, 2]. The concept has been extended to multivariate Granger Causality approaches, often combined with a sparsity inducing Lasso penalty, cf. [1, 33].

Beyond Granger Causality, there have been many recent approaches to Causal Discovery for time series, particularly at the FG level. PCMCI/PCMCI+ [26, 29, 27] and the related LPCMCI [8] execute a two-step procedure. The first step consists of estimating a set of parents for each variable, which is based on PC and FCI, respectively. In the second step, we test for conditional independence of any two variables conditioned on the union of their parents. Furthermore, VAR-LiNGAM [13], DYNOTEARS [22] and SVAR-(G)FCI [17] are vector-autoregressive extensions of LiNGAM, NoTEARS and FCI / GFCI, respectively. Further recent references for time series causal discovery methods can be found in [24, 6, 12, 38].

2.2 Structural causal models

Next we introduce *Structural Causal Models* (SCM) [23], also known as *Structural Equation Models* or *Functional Causal Models*. These models assume that child nodes in a causal graph have a functional dependence on their parents. More precisely, given a set of variables X_1, \dots, X_m , each variable X_i can be represented in terms of some function F_i and its parents $\mathcal{P}(X_i)$ as

$$X_i = F_i(\mathcal{P}(X_i), N_i), \quad (1)$$

where N_i are independent noise terms with a given distribution. In practice, the SCM in Eq. (1) is often too general and one considers a more restricted class. In particular, we focus on *Causal Additive Models* (CAM) [3], where both F_i and the noise are additive, such that (for univariate functions f_{ij})

$$X_i = \sum_{X_j \in \mathcal{P}(X_i)} f_{ij}(X_j) + N_i. \quad (2)$$

A special case are linear causal models with $f_{ij}(x) = \beta_{ij}x$, s.t. $X_i = \sum_{X_j \in \mathcal{P}(X_i)} \beta_{ij}X_j + N_i$.

In the case of time series, the functions F_i can in principle depend on time, which creates additional difficulties for estimating causal relationships between variables. Within the proposed framework, we will assume that the functions F_i are invariant over time and that, in particular, the causal dependence between variables does not change over time:

$$X_i(t) = \sum_{X_j(t') \in \mathcal{P}(X_i(t))} f_{ij}(X_j(t')) + N_i, \quad (3)$$

where necessarily $t' \leq t$ for each $X_j(t') \in \mathcal{P}(X_i(t))$, in order to preserve the order of time.

In general, one cannot fully resolve the causal graph from observational data generated by a fully general SCM as in Eq. (1). However, if the model is restricted to specific classes, such as a non-linear CAM, the full DAG is, in principle, identifiable [3]. Naturally, in a real use case the class of models must be assumed or known *a priori*. Moreover, while a linear SCM with Gaussian noise renders the DAG unidentifiable for IID data, this is not always the case for time series [25, Chapter 10].

2.3 Related work

When novel methods are developed, the performance is usually evaluated on synthetic data (cf. [24]), which helps to identify strengths and weaknesses of the methods. Unfortunately, the results cannot be directly compared to other methods, as the generated data is often not made available. For this reason, it is desired to create general benchmarks that can be used to compare methods. One example is the benchmark created for the ChaLearn challenge for pairwise causal discovery [11]. Recent work has been done to create a unified benchmark for causal discovery for time series. CauseMe [28] contains a mix of synthetic, hybrid, and real data based on challenges found in climate and weather data. For these scenarios, the platform provides a ground truth (based on domain knowledge for real data).

We see three points to how our proposed framework goes beyond the capabilities of CauseMe. First, CauseMe is based on a “static” set of data used for benchmarking results. This increases the chance of “overfitting” new methods to perform well on the specific use case covered in the benchmark, rather than to perform well in general. Our proposed framework allows one to generate vast amounts of data with different properties, including number of observations and number of variables, enabling the user to select more robust hyperparameters that perform well under a diverse selection of problems. Second, our framework provides a greater flexibility to the user, which allows them to understand the behavior of the method in specific edge cases (e.g., when underlying assumptions are violated),

or how a method scales with the number of time series. Third, the proposed framework contributes to reproducibility. It allows the user to specify the configuration used for an experiment, based on which others can regenerate the very same data, facilitating the reproduction of their results.

3 Data generating process

We now describe the proposed data generating process. The general idea follows three steps: (i) specify and generate a time series causal graph, (ii) specify and generate a structural causal model (SCM), and (iii) specify the noise and runtime configuration to generate synthetic time series data. We expose a hierarchical `DataGenerationConfig` object, which contains subconfiguration objects for the causal graph (`CausalGraphConfig`), for the SCM (`FunctionConfig`), and to generate data (`NoiseConfig` and `RuntimeConfig`). See Appendix C for a complete example.

Figure 1 captures a high-level overview of the process to go from a partially defined configuration to generated data, while Algorithm 1, Algorithm 2, and Algorithm 3 in Appendix A detail the full process. The concept of *complexity* exists in the configurations. Each configuration object allows the user to make the problem as complex or simple as they would like. Using the `complexity` parameter allows the system to specify default values if the user does not want to fully specify the configuration.

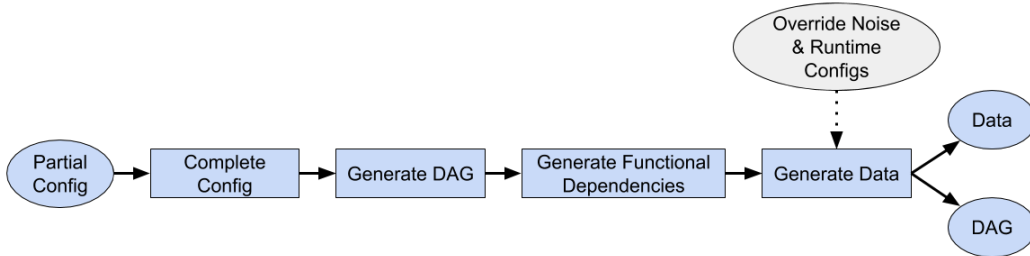


Figure 1: The user provides a `DataGenerationConfig` object. This can be partially defined and default values based on the complexity setting will be used to complete the configuration. Given the completed configuration, a time series causal graph (a Full Time Graph, cf. § 2.1) is randomly generated. For each edge of the DAG, a functional dependency is randomly chosen, resulting in a randomly generated SCM from which data can be generated. Multiple data sets with varying number of observations can be returned for a single SCM. Therefore, the user receives a list of data sets and a single DAG. Optionally, the user can override the original `NoiseConfig` and/or `RuntimeConfig` provided with the `DataGenerationConfig`. E.g., the user can change the noise distributions or regenerate a data set to also return unobserved variables.

We expose a configuration for four types of variables: targets, features, latent, and noise. This provides the user with the capability to define the structure around specific variables. For example, when performing causal feature selection, such as in SyPI [19], one may want to ensure a target variable is a sink node, i.e., it has no children. This is a key feature of the graph, function, and noise configurations as they allow one to fully specify the model based on the assumptions of ones’ method, such as causal sufficiency and linearity for DYNOTEARS [22]. Additionally, sparsity of the causal graph can easily be controlled by specifying the likelihood of edges (as a universal setting or for each variable type) and the maximum number of parents and children for each variable type.

Figure 2 (a) shows causal sufficiency being broken as we have introduced an unobserved variable. Figure 2 (b) displays the time series for the three observed variables. However, if desired, it is possible to return all the synthetic data, including the latent variables and the noise variables by setting `return_observed_data_only` to `False` in the `RuntimeConfig`. Another powerful feature is defining the noise distributions after creating the SCM. The user can easily generate new data with different noise distributions and/or signal-to-noise ratio without needing to create a new SCM as depicted in Figure 1. Finally, the process is open source and an example script is provided to demonstrate the effects of the various configuration settings, to provide further information on the complexity settings, and to allow users to generate data for their use cases.¹

¹<https://github.com/causalens/cdml-neurips2020>

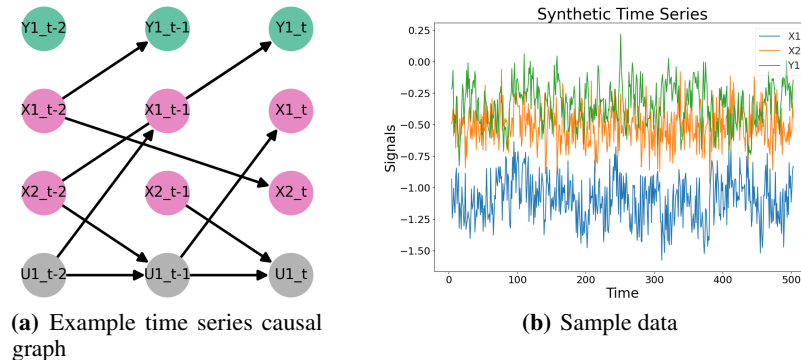


Figure 2: (a) A simple time series causal graph (a Full Time Graph, cf. § 2.1) with a maximum lag of two, one target variable (Y_1), two feature variables (X_1 and X_2), and one latent variable (U_1); the noise nodes are not displayed for simplicity. U_1 shows an autoregressive relationship. (b) The three *observed* time series drawn from a SCM with the causal graph in (a).

4 Experiments

In order to demonstrate the breadth of data that can be created using the process proposed in § 3, we performed several experiments (see Table 1 below and Table 4 in Appendix B), each of which invalidates a specific assumption of the methods under test. The goal is not to fully evaluate each method, but to show how the proposed data generating process supports testing.

The established methods were chosen due to their popularity and the availability of open source Python implementations. Additionally, we evaluated a custom implementation of a multivariate version of Granger Causality [1, 33]. We also tested a bivariate version of Granger Causality [10] and the PC [34] algorithm (based on PCMCI [29]). We do not report results for these two methods, as they significantly underperformed. Supplemental experiments and results are provided in Appendix B.

The methods and hyperparameters are listed in Table 2. PCMCI(+) was used with Partial Correlation and a threshold of $\alpha = 0.02$ for statistical significance. Note that the hyperparameters have been chosen manually to avoid an overly “aggressive” selection of links, which would result in high FPRs or FNRs (cf. § 4.1). We note that the results are sensitive to hyperparameter choices, generally resulting in a tradeoff between higher TPRs and TNRs (cf. Figure 6 and Figure 7).

For each experiment, 200 unique SCMs were generated from the same parameterization space defined for the specific experiment and a single data set with 1000 samples was generated from each SCM. For the causal sufficiency experiment, the number of feature nodes is kept at 10, while the number of latent nodes increases from 0 to 20. For the other experiments, the parameterization space allows for a variable number of nodes to not limit the experiments to a single graph size.² Additionally, the metrics are normalized to allow for a fair comparison between varying graph sizes. The results presented in § 4.2 capture the average metrics (with a maximum lag to consider of $\ell_{\max} = 5$), defined in § 4.1, for each causal discovery method across the 200 data sets with known causal ground truth. The synthetic data never contained true lags longer than 5 timesteps in the past. Finally, to demonstrate the effect of the choice of hyperparameters on PCMCI and DYNOTEARS, we performed the causal sufficiency experiment using 100 unique SCMs and a single data set with 500 samples generated from each SCM while modifying two hyperparameters of each method.

Table 1: Experiments

Name	Description
1. Causal Sufficiency	The number of observed variables remains fixed while the number of latent variables is increased.
2. Non-Linear Dependencies	The likelihood of linear functions is decreased while the likelihood of monotonic and periodic functions is increased.
3. Instantaneous Effects	The minimum allowed lag in the SCM is reduced from 1 to 0.

²Experiment data sets are provided: <https://github.com/causalens/cdml-neurips2020>

Table 2: Causal discovery methods and their chosen hyperparameters

Name	Source	Hyperparameters
PCMCI [29]	Python package tigramite version 4.2	tau_min = 0, tau_max = 5, pc_alpha = 0.01
PCMCI+ [27]	Python package tigramite version 4.2	tau_min = 0, tau_max = 5, pc_alpha = 0.05
DYNOTEARS [22]	Python package causalnex	lambda_w = lambda_a = 0.15, w_threshold = 0.0, p = 5
VAR-LiNGAM [13]	Python package lingam	prune = True, criterion = 'aic', lags = 5
Multivariate Granger Causality [1, 33]	Cross-validated Lasso regression (scikit-learn) + one-sided t-test (statsmodels)	cv_alphas = [0.02, 0.05, 0.1, 0.2, 0.3, 0.4, 0.5], max_lag = 5

4.1 Metrics and evaluation details

In order to evaluate a time series causal graph, we must define a maximum lag variable ℓ_{\max} , which controls the largest potential lag we want to evaluate. Given m time series X_1, \dots, X_m , we define the “universe” of variables

$$\mathbb{X}_t := \left\{ X_i(t-s) \mid i = 1, \dots, m \text{ and } s = 0, \dots, \ell_{\max} \right\} \quad (4)$$

and the set of possible links with maximal lag ℓ_{\max} is given by

$$\mathbb{L}_t := \left\{ X_i(t-s) \rightarrow X_j(t) \mid i, j \in \{1, \dots, m\} \text{ and } s = 0, \dots, \ell_{\max} \text{ s.t. } s > 0 \text{ or } i \neq j \right\}. \quad (5)$$

As a special case, Eq. (5) contains the instantaneous links

$$\tilde{\mathbb{L}}_t := \{ X_i(t) \rightarrow X_j(t) \mid i, j \in \{1, \dots, m\}, i \neq j \}. \quad (6)$$

Note that the latter coincides with all the possible links $m(m-1)$ prominent in an IID setup. The total number of links is given by $|\mathbb{L}_t| = (\ell_{\max} + 1)m^2 - m$. Due to the acyclicity constraint for instantaneous links $\tilde{\mathbb{L}}_t$, a valid DAG can contain at most $\ell_{\max}m^2 + m(m-1)/2$ links.

One can then define the True Positives (TP) as the correctly identified links in \mathbb{L}_t , and similarly for True Negatives (TN), False Positives (FP) and False Negatives (FN). Table 3 contains the metrics used to evaluate the performance, including NTP, NFP and NFN, which can be used to derive SHD, F1, as well as the True Positive Rate: $\text{TPR} = \text{NTP} / (\text{NTP} + \text{NFN})$. Using these normalized values allows for a more intuitive understanding of how each component of the SHD and F1 metrics behave. Note that $\text{SHD} = \text{NFP} + \text{NFN}$, such that a SHD of 5% means that the method misclassified 5% of the edges. For a typical graph we have much more Negatives (N) than Positives (P), such that the NTP and NFN will be on a much smaller scale than TPR and FNR, making a direct comparison impossible. However, NFP and False Positive Rate (FPR) should be roughly on the same scale.

Table 3: Metrics

Name	Acronym	Description
F1-Score	F1	Calculated as $F1 = \text{TP} / (\text{TP} + (\text{FP} + \text{FN}) / 2)$.
Structural Hamming Distance [36]	SHD	Normalized as $(\text{FP} + \text{FN}) / \mathbb{L}_t $.
Normalized True Positives	NTP	Calculated as $\text{TP} / \mathbb{L}_t $.
Normalized False Positives	NFP	Calculated as $\text{FP} / \mathbb{L}_t $.
Normalized False Negatives	NFN	Calculated as $\text{FN} / \mathbb{L}_t $.

4.2 Results

Performance of the causal discovery methods under evaluation against the metrics defined in § 4.1 are provided in Figure 3, Figure 4, and Figure 5 for the experiments defined in Table 1, respectively. Figure 3 and Figure 4 show how all the methods are affected by latent variables and non-linear dependencies, respectively. The compared methods share the assumptions of causal sufficiency and linearity, and, as expected, their performance degrades with the violation of these assumptions. The introduction of non-linearities decreases the methods’ performance to a much larger extent than the invalidation of causal sufficiency. However, the choice of hyperparameters has as much of an effect, even when assumptions are met. Not all hyperparameters have the same importance and, accordingly, we suggest a range of values in Figure 6 and Figure 7 for PCMCI and DYNOTEARS, respectively. Finally, results for the supplemental experiments are provided in Appendix B.

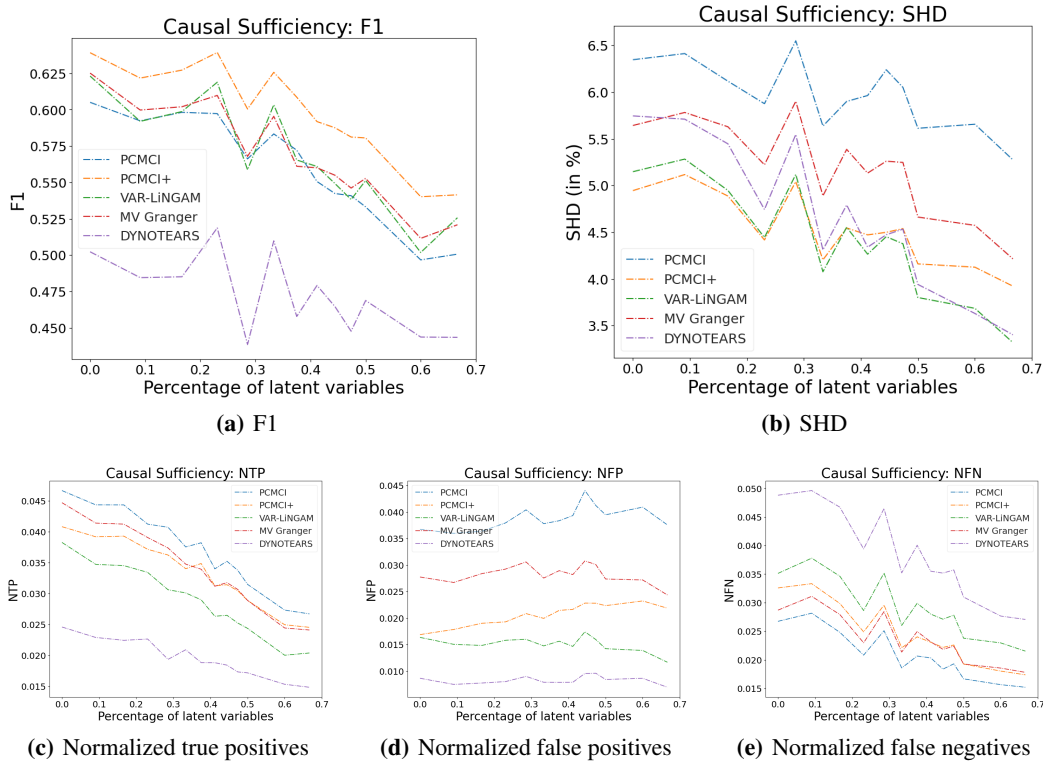


Figure 3: Causal Sufficiency - The percentage of latent variables is the number of latents divided by the number of observed plus latent. (a) F1 decreases for all methods as more latent variables are added. (b) Note that SHD also decreases. As defined in Table 3, SHD is only a function of FPs and FNs, while F1 is also a function of TPs. (c) and (e) respectively show that TPs and FNs decrease at a similar rate as more latent variables are added. TPs have a larger effect on F1, hence why we observe an overall decrease. (d) The relative minor changes in FPs when compared to the larger decrease in FNs leads to an overall decrease in SHD.

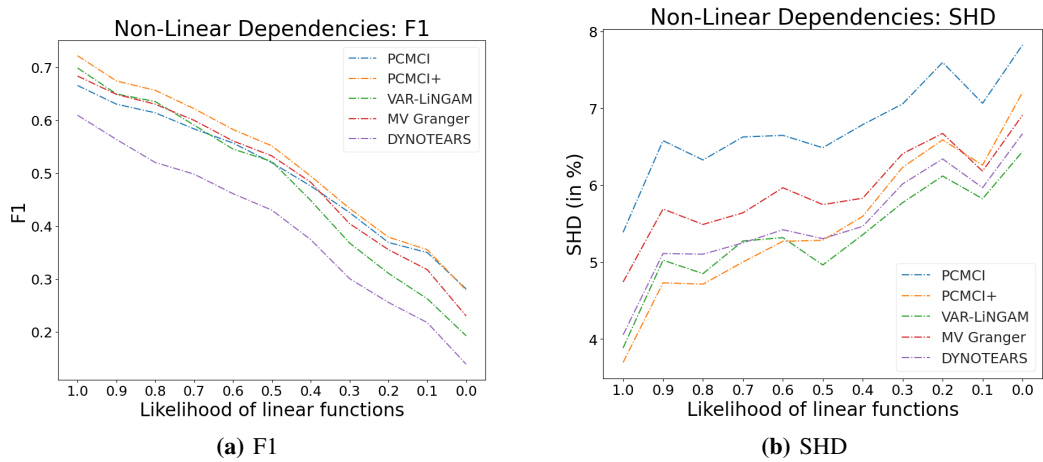


Figure 4: Non-Linear (Monotonic and Periodic) Dependencies - (a) F1 decreases and (b) SHD increases as the percentage of linear dependencies in the system is decreased. A key observation is that VAR-LiNGAM, multivariate Granger, and DYNOTEARS have the largest decreases in F1, which is expected as they are linear vector autoregressive models. PCMCI and PCMCI+ are negatively impacted by the use of Pearson correlation within the conditional independence test. We explicitly apply linear methods to a non-linear setting to demonstrate how the methods perform when this assumption is violated; we are not expecting linear methods to perform well for non-linear problems.

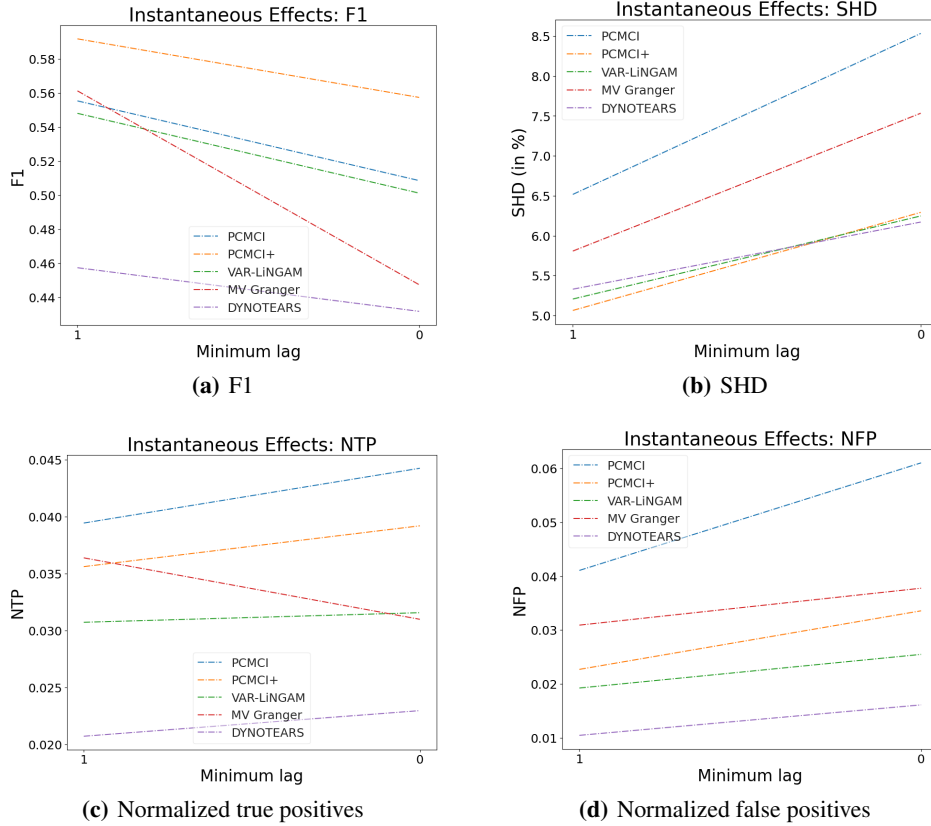


Figure 5: Instantaneous Effects - (a) F1 decreases and (b) SHD increases for all methods when instantaneous effects are allowed. By construction, multivariate Granger does not regress on covariates with lag 0 and hence cannot identify instantaneous effects; the TPs (c) drop substantially. PCMCI+ is specifically designed to handle instantaneous links [27]. PCMCI is capable of returning edges at lag 0, but the edges are not oriented. As such, for each TP, there will be a FP, as seen by the steep increase in FPs (d) for PCMCI when compared to the slower increase for PCMCI+. Note that this is one scenario when PCMCI does not return a valid DAG, but a CPDAG (cf. § 2).

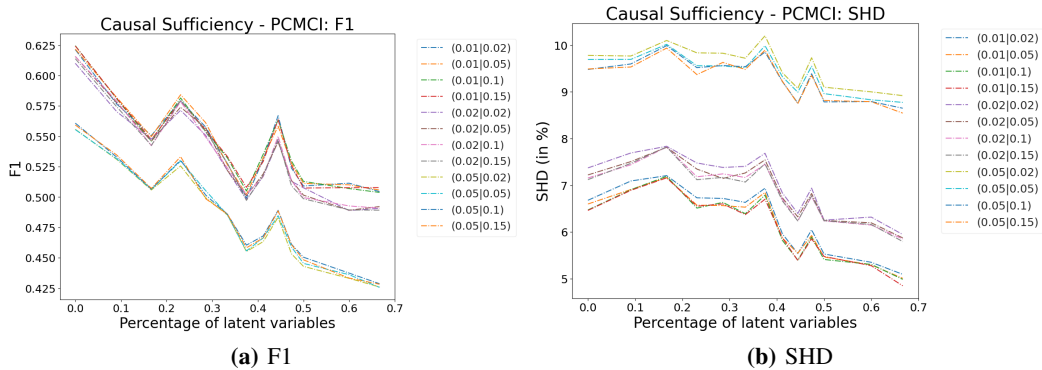


Figure 6: (a) F1 and (b) SHD for different choices of hyperparameters for PCMCI. We varied the p-value threshold for the final selection (first parameter in the legend) over $\{0.01, 0.02, 0.05\}$ and the p-value for the PC algorithm (second parameter in the legend) over $\{0.02, 0.05, 0.1, 0.15\}$. Both F1 and SHD show that the variation of the p-value for the PC step only has a minor effect on the performance, and the majority is accounted for by the first parameter, which is the p-value threshold used for the final score. From the results, the best performance is achieved for the final p-value threshold equal to 0.01 and the p-value for the PC step in the range of $[0.02, 0.2]$.

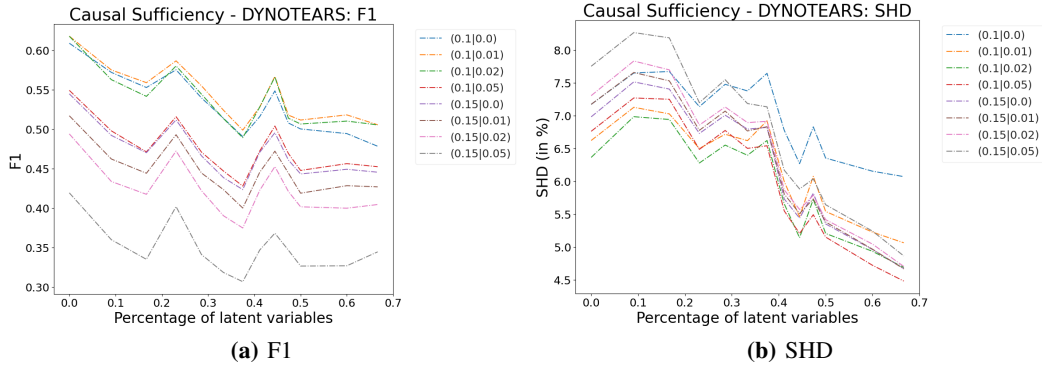


Figure 7: (a) F1 and (b) SHD for different choices of hyperparameters for DYNOTEARS. The first parameter in the legend is equal to the L1-penalty on all the variables (specified by λ_w and λ_a , which are here chosen to be equal). The second parameter is $w_{\text{threshold}}$, which is used as a threshold for the cutoff of weights. The results suggest that the optimal choices are moderate values of the L1-penalties (λ_w and λ_a between 0.5 and 1.0), combined with a small value for $w_{\text{threshold}}$.

5 Conclusion

We have proposed a process to generate synthetic time series data with a known ground truth causal structure and demonstrated its functionality by evaluating prominent causal discovery techniques for time series data. The process is easily parameterizable yet provides the capability to generate data from vastly different scenarios. The process is open source and an example script to allow users to generate data for their use cases is provided. This is the main contribution over existing benchmarks, such as CauseMe [28], as their scope is restricted with respect to number of observations, number of variables, and specific dynamic system challenges. We have demonstrated how the proposed framework can be used to discriminate the performance of different causal discovery methods under a variety of conditions and how this framework can be used for fine tuning hyperparameters without the fear of overfitting to a target benchmark.

We believe our proposed data generating process can be used to support both research and the practical applications of causal discovery. As a result of our experiments, we have identified two important research directions: (i) the continued development of efficient, non-linear methods, and (ii) less reliance on hyperparameters. To address sparsity of data and lack of ground truth, a researcher/practitioner will generate synthetic data in agreement with their domain knowledge. The resulting synthetic benchmark will provide a principled approach to the development/selection of a method and its hyperparameters, as opposed to simply overfitting to one’s data.

Future work The current implementation can be extended in various directions, particularly around functional forms, noise, and dynamic graphs. Currently, the process only supports additive, homoscedastic noise; adding support for multiplicative and heteroskedastic noise would be beneficial. The current process also produces static models. Allowing for the distributions, function parameters, and causal graph to change over time will produce synthetic data with changepoints, regime shifts, and/or interventions.

Acknowledgments and Disclosure of Funding

The authors would like to thank Microsoft for Startups for supporting this research through their contribution of GitHub and Azure credits. We would also like to thank the two anonymous reviewers whose comments helped us to improve and clarify the paper.

References

- [1] Andrew Arnold, Yan Liu, and Naoki Abe. Temporal causal modeling with graphical granger methods. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 66–75, 2007.
- [2] Lionel Barnett, Adam B Barrett, and Anil K Seth. Granger causality and transfer entropy are equivalent for gaussian variables. *Physical review letters*, 103(23):238701, 2009.
- [3] Peter Bühlmann, Jonas Peters, Jan Ernest, et al. CAM: Causal additive models, high-dimensional order search and penalized regression. *The Annals of Statistics*, 42(6):2526–2556, 2014.
- [4] David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of machine learning research*, 3(Nov):507–554, 2002.
- [5] Diego Colombo, Marloes H Maathuis, Markus Kalisch, and Thomas S Richardson. Learning high-dimensional directed acyclic graphs with latent and selection variables. *The Annals of Statistics*, pages 294–321, 2012.
- [6] Doris Entner and Patrik O Hoyer. On causal discovery from time series data using fci. *Probabilistic graphical models*, pages 121–128, 2010.
- [7] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*, volume 1. Springer series in statistics New York, 2001.
- [8] Andreas Gerhardus and Jakob Runge. High-recall causal discovery for autocorrelated time series with latent confounders. *arXiv preprint arXiv:2007.01884*, 2020.
- [9] Clark Glymour, Kun Zhang, and Peter Spirtes. Review of causal discovery methods based on graphical models. *Frontiers in genetics*, 10:524, 2019.
- [10] Clive WJ Granger. Investigating causal relations by econometric models and cross-spectral methods. *Econometrica: journal of the Econometric Society*, pages 424–438, 1969.
- [11] Isabelle Guyon, Alexander Statnikov, and Berna Bakir Batu. *Cause Effect Pairs in Machine Learning*. Springer, 2019.
- [12] Biwei Huang, Kun Zhang, Jiji Zhang, Joseph Ramsey, Ruben Sanchez-Romero, Clark Glymour, and Bernhard Schölkopf. Causal discovery from heterogeneous/nonstationary data. *Journal of Machine Learning Research*, 21(89):1–53, 2020.
- [13] Aapo Hyvärinen, Kun Zhang, Shohei Shimizu, and Patrik O Hoyer. Estimation of a structural vector autoregression model using non-gaussianity. *Journal of Machine Learning Research*, 11(5), 2010.
- [14] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [15] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [16] Yann LeCun, Corinna Cortes, and CJ Burges. MNIST handwritten digit database. *ATT Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- [17] Daniel Malinsky and Peter Spirtes. Learning the structure of a nonstationary vector autoregression. *Proceedings of machine learning research*, 89:2986, 2019.
- [18] Daniele Marinazzo, Mario Pellicoro, and Sebastiano Stramaglia. Kernel-granger causality and the analysis of dynamical networks. *Physical review E*, 77(5):056215, 2008.
- [19] Atalanti A Mastakouri, Bernhard Schölkopf, and Dominik Janzing. Necessary and sufficient conditions for causal feature selection in time series with latent common causes. *arXiv preprint arXiv:2005.08543*, 2020.
- [20] Christopher Meek. *Graphical Models: Selecting causal and statistical models*. PhD thesis, PhD thesis, Carnegie Mellon University, 1997.
- [21] Ignavier Ng, AmirEmad Ghassami, and Kun Zhang. On the role of sparsity and dag constraints for learning linear dags. *arXiv preprint arXiv:2006.10201*, 2020.
- [22] Roxana Pamfil, Nisara Sriwattanaworachai, Shaan Desai, Philip Pilgerstorfer, Konstantinos Georgatzis, Paul Beaumont, and Bryon Aragam. Dynotears: Structure learning from time-series data. In *International Conference on Artificial Intelligence and Statistics*, pages 1595–1605, 2020.

- [23] Judea Pearl et al. Causal inference in statistics: An overview. *Statistics surveys*, 3:96–146, 2009.
- [24] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. Causal inference on time series using restricted structural equation models. In *Advances in Neural Information Processing Systems*, pages 154–162, 2013.
- [25] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference*. The MIT Press, 2017.
- [26] Jakob Runge. Causal network reconstruction from time series: From theoretical assumptions to practical estimation. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 28(7):075310, 2018.
- [27] Jakob Runge. Discovering contemporaneous and lagged causal relations in autocorrelated nonlinear time series datasets. In *Proceedings of the Thirty-Sixth Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 1388–1397. AUAI Press, 03–06 Aug 2020.
- [28] Jakob Runge, Sebastian Bathiany, Erik Bollt, Gustau Camps-Valls, Dim Coumou, Ethan Deyle, Clark Glymour, Marlene Kretschmer, Miguel D Mahecha, Jordi Muñoz-Marí, et al. Inferring causation from time series in earth system sciences. *Nature communications*, 10(1):1–13, 2019.
- [29] Jakob Runge, Peer Nowack, Marlene Kretschmer, Seth Flaxman, and Dino Sejdinovic. Detecting and quantifying causal associations in large nonlinear time series datasets. *Science Advances*, 5(11), 2019.
- [30] Thomas Schreiber. Measuring information transfer. *Physical review letters*, 85(2):461, 2000.
- [31] Shohei Shimizu, Patrik O Hoyer, Aapo Hyvärinen, and Antti Kerminen. A linear non-gaussian acyclic model for causal discovery. *Journal of Machine Learning Research*, 7(Oct):2003–2030, 2006.
- [32] Shohei Shimizu, Takanori Inazumi, Yasuhiro Sogawa, Aapo Hyvärinen, Yoshinobu Kawahara, Takashi Washio, Patrik O Hoyer, and Kenneth Bollen. Directlingam: A direct method for learning a linear non-gaussian structural equation model. *The Journal of Machine Learning Research*, 12:1225–1248, 2011.
- [33] Ali Shojaie and George Michailidis. Discovering graphical granger causality using the truncating lasso penalty. *Bioinformatics*, 26(18):i517–i523, 2010.
- [34] Pater Spirtes, Clark Glymour, Richard Scheines, Stuart Kauffman, Valerio Aimale, and Frank Wimberly. Constructing bayesian network models of gene expression networks from microarray data. 2000.
- [35] Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. *Causation, prediction, and search*. MIT press, 2000.
- [36] Ioannis Tsamardinos, Laura E Brown, and Constantin F Aliferis. The max-min hill-climbing bayesian network structure learning algorithm. *Machine learning*, 65(1):31–78, 2006.
- [37] Gherardo Varando. Learning dags without imposing acyclicity. *arXiv preprint arXiv:2006.03005*, 2020.
- [38] Sebastian Weichwald, Martin E. Jakobsen, Phillip B. Mogensen, Lasse Petersen, Nikolaj Thams, and Gherardo Varando. Causal structure learning from time series: Large regression coefficients may predict causal links better in practice than small p-values. volume 123 of *Proceedings of the NeurIPS 2019 Competition and Demonstration Track, Proceedings of Machine Learning Research*, pages 27–36. PMLR, 2020.
- [39] Yue Yu, Jie Chen, Tian Gao, and Mo Yu. Dag-gnn: Dag structure learning with graph neural networks. *arXiv preprint arXiv:1904.10098*, 2019.
- [40] Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. Dags with no tears: Continuous optimization for structure learning. In *Advances in Neural Information Processing Systems*, pages 9472–9483, 2018.
- [41] Xun Zheng, Chen Dan, Bryon Aragam, Pradeep Ravikumar, and Eric Xing. Learning sparse nonparametric dags. In *International Conference on Artificial Intelligence and Statistics*, pages 3414–3425. PMLR, 2020.

A Algorithmic representation of data generating process

Algorithm 1, Algorithm 2, and Algorithm 3 define the data generating process proposed in § 3 and provide the low-level steps to go from a configuration to generated data as shown in Figure 1.

Algorithm 1: Time Series Data Generation

Input: config: DataGenerationConfig

```

1 Complete missing values with defaults for config.noise_config (based on complexity value)
  and config.runtime_config.
2 Generate TimeSeriesCausalGraph per Algorithm 2.
3 Generate StructuralCausalModel per Algorithm 3.
4 foreach num_samples and data_generating_seed in config.runtime_config do
5   Seed process using data_generating_seed.
6   Initialize all data with zeroes.
7   foreach Noise variable  $N_i$  in StructuralCausalModel do
8     if noise_config.noise_variance is provided as a range then
9       | noise_var  $\sim$  Uniform(noise_config.noise_variance)
10    else
11      | noise_var  $\leftarrow$  noise_config.noise_variance
12    end
13    Randomly sample noise distribution from noise_config.distributions with
      probabilities defined in noise_config.prob_distributions.
14     $N_i \leftarrow$  Sample num_samples IID samples from the chosen distribution with variance
      noise_var.
15  end
16  foreach Noise variable  $N_i$  with an autoregressive edge in StructuralCausalModel do
17    for  $t \leftarrow 1$  to num_samples do
18      |  $N_i(t) \leftarrow N_i(t) + f_i(N_i(t-1))$ 
19    end
20  end
21  for  $t \leftarrow$  config.graph_config.max_lag to num_samples do
22    foreach Non-noise variable  $X_i$  in topological order of graph do
23      | foreach Parent variable  $X_j$ , lagged time index  $t'$ , and functional dependency  $f_{ij}$  do
24        |  $X_i(t) \leftarrow X_i(t) + f_{ij}(X_j(t'))$ 
25        | Note: This includes the additive noise as the noise is a parent and its functional
          dependency is the identity function as described in Algorithm 3.
26      | end
27    end
28  end
29 end

```

Output: Tuple[List[Dataset], TimeSeriesCausalGraph]

Algorithm 2: Time Series Causal Graph Generation

Input: graph_config: CausalGraphConfig

```
1 Complete missing values with defaults based on complexity value for graph_config.
2 Initialize empty DAG with  $M$  nodes, where  $M = (1 + \text{graph\_config.include\_noise}) *
  (\text{graph\_config.num\_targets} + \text{graph\_config.num\_features} +
  \text{graph\_config.num\_latent}) * (1 + \text{graph\_config.max\_lag})$ .
3 if graph_config.include_noise then Add connections for noise.
4   foreach noise_node do
5     | Add edge to its corresponding target, feature, or latent node.
6   end
7 end
8 if graph_config.max_lag > 0 then Add autoregressive edges.
9   foreach target, feature, and noise variable in graph do
10    | add_edge  $\sim$  Bernoulli(graph_config.prob_<var_type>_autoregressive)
11    | if add_edge is True then
12    | | Add forward (in time) edge between consecutive nodes of current variable.
13    | end
14  end
15 end
16 foreach Non-noise node at time  $t$  (i.e., lag of 0) do
17   | Construct a list of possible parents based on graph_config.min_lag (or current topology
18   | of DAG if min_lag is 0 to maintain acyclic graph),
19   | graph_config.allow_latent_direct_target_cause, and
20   | graph_config.allow_target_direct_target_cause.
21   | Shuffle list of possible parents.
22   | while Current number of parents < graph_config.max_<var_type>_parents and length
23   | of list of possible parents > 0 do
24     | Pop first element from list of possible parents.
25     | prob_edge  $\leftarrow$  graph_config.prob_<var_type>_parent if not None else
26     | graph_config.prob_edge
27     | add_edge  $\sim$  Bernoulli(prob_edge) # Note: prob_edge controls graph sparsity.
28     | if add_edge is True and the number of existing children of the current possible parent
29     | node < graph_config.max_<parent_type>_children then
30       | Add edge from possible parent to current node.
31       | Add edges between nodes with appropriate lags up to graph_config.max_lag,
32       | e.g., if  $X_j(t-1) \rightarrow X_i(t)$ , then  $X_j(t-2) \rightarrow X_i(t-1)$ , etc.
33       | if Number of parents of current node from the same variable  $\geq$ 
34       | graph_config.max_parents_per_variable then
35       | | Remove any other instances of the same variable as the current parent from the
36       | | list of possible parents. For example, this parameters prevents
37       | |  $X_j(t-1) \rightarrow X_i(t)$  and  $X_j(t-2) \rightarrow X_i(t)$  if
38       | | graph_config.max_parents_per_variable is 1.
39       | end
40     | end
41   | end
42 end
43 end
44 Output: TimeSeriesCausalGraph
```

Algorithm 3: Structural Causal Model Generation

Input: function_config: FunctionConfig, causal_graph: TimeSeriesCausalGraph

```
1 Complete missing values with defaults based on complexity value for function_config.
2 foreach Node in causal_graph at time t (i.e., lag of 0) do
3   if Current node is a noise node and it has a parent then
4     Sample linear weights for functional dependency  $f_i$  between its parent and itself as its
       only possible parent is the lagged version of itself and this relationship is currently
       limited to linear.
5     Set autoregressive relationship to linear function with sampled parameters.
6   else
7     foreach Parent of current node do
8       if Current parent is a noise node then
9         Set functional dependency  $f_{ij}$  to the identity function as noise is currently simply
           treated as additive.
10      else
11        Randomly sample function type from function_config.functions with
           probabilities function_config.prob_functions.
12        Randomly sample parameters for chosen function type (cf. example script for
           details).
13        Set functional dependency  $f_{ij}$  to sampled function with sampled parameters.
14      end
15    end
16  end
17 end
```

Output: StructuralCausalModel

B Supplemental experiments

We performed two additional experiments, as outlined in Table 4, on the methods in Table 2. As described in § 4, for each experiment, 200 unique SCMs were generated from the same parameterization space defined for the specific experiment. For the non-Gaussian noise experiment, a single data set with 1000 samples was generated from each SCM, while we varied the number of samples in the data set from each SCM for the IID experiment. The following results shown in Figure 8 and Figure 9 capture the average metrics, defined in § 4.1, for each causal discovery method across the 200 data sets with known causal ground truth for the experiments defined in Table 4, respectively.

Table 4: Supplemental experiments

Name	Description
4. IID Data	Only IID data is generated.
5. Non-Gaussian Noise	The likelihood of additive Gaussian noise is decreased while the likelihood of Laplace, Student’s t, and uniform distributed noise is increased. The noise variance remains unchanged.

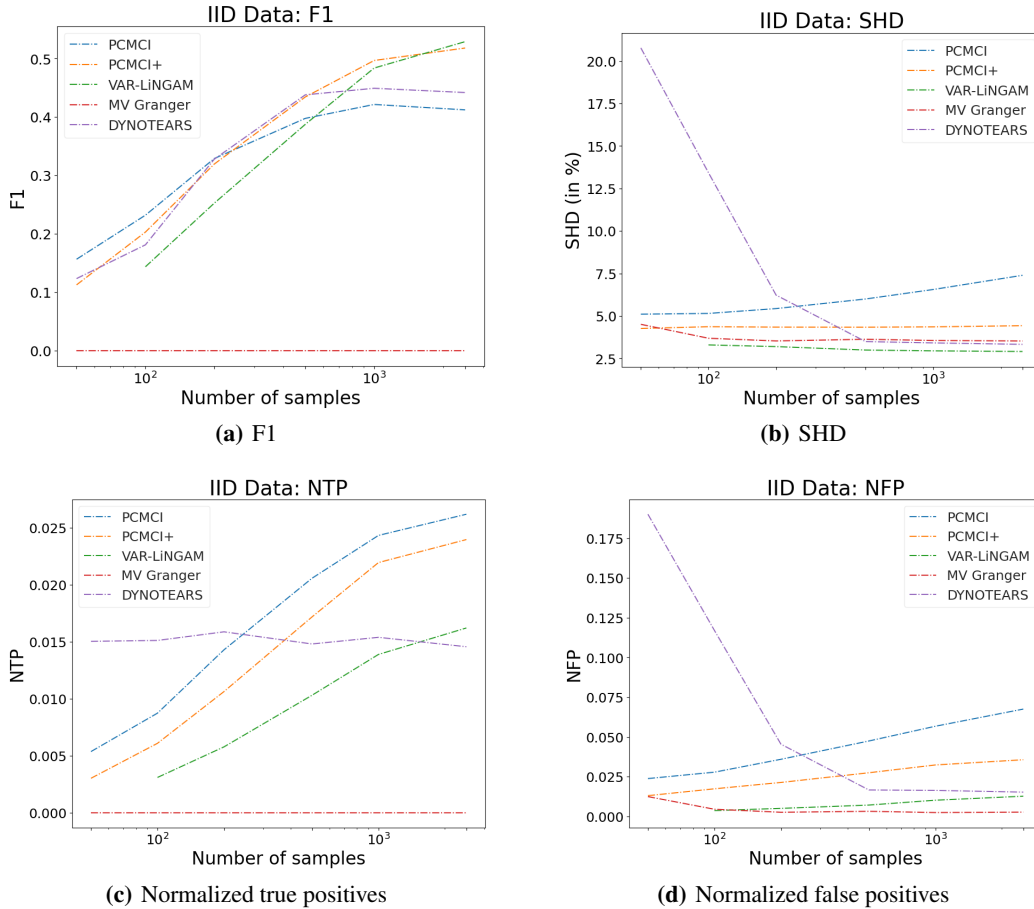


Figure 8: IID Data - For IID data the true links in the causal graph only exist between nodes with a relative lag of 0. As described in Figure 5, our implementation of multivariate Granger can never identify instantaneous effects. For the IID case, all the true edges are instantaneous effects and the method returns zero true positives (c). (a) F1 improves for all other methods as the number of observations increase. The increase in SHD (b) for PCMCI is again attributed to the increase in false positives (d) due to its inability to orient edges of contemporaneous links as mentioned in Figure 5.

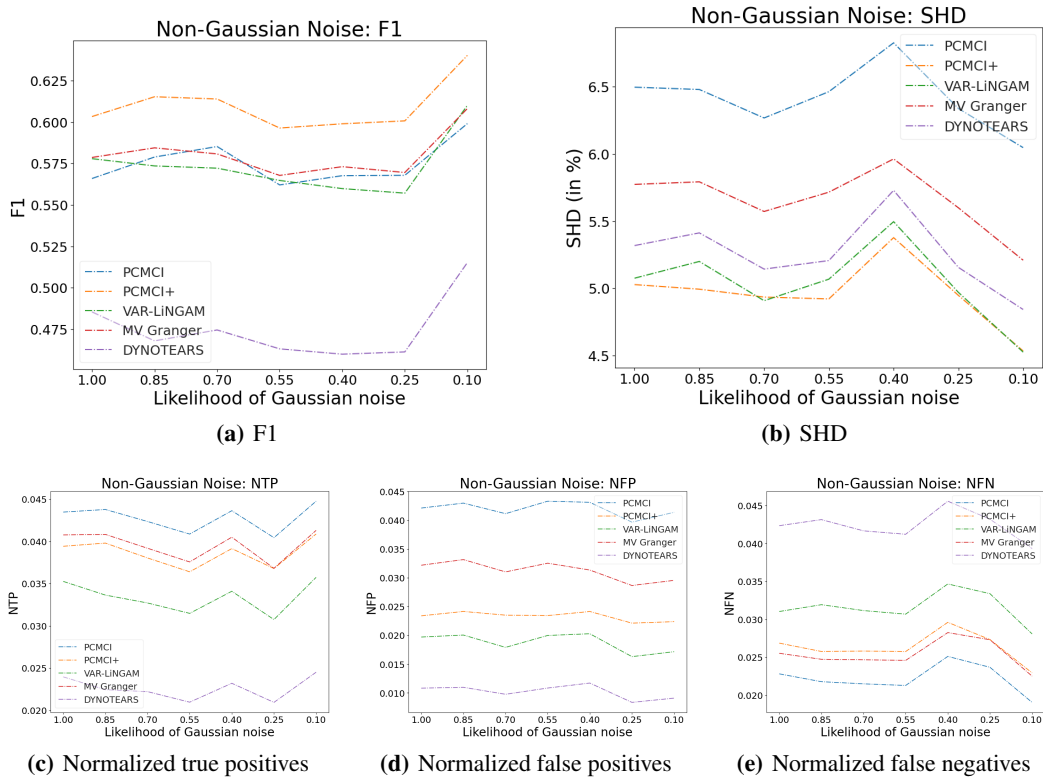


Figure 9: Non-Gaussian Noise - (a) F1 is stable for all methods so they perform fairly well when noise is sampled from distributions with fatter tails. The increase in F1 below 25% can be attributed to both the increase in true positives (c) and the decrease in false negatives (e). (b) SHD follows the trajectory of the false negatives as the false positives (d) are relatively flat in comparison.

C Example configuration

```
import (CausalGraphConfig, DataGenerationConfig, FunctionConfig, NoiseConfig,
        RuntimeConfig)

# DataGenerationConfig is the top-level configuration object and provides the capability
# to specify all necessary parameters to generate a SCM and sample time series data.
config = DataGenerationConfig(
    # Controls random behavior and ensures reproducibility for graph and SCM generation.
    random_seed=1,
    # Used to initialize any unspecified configurations.
    # They are all initialized in this example so value would be ignored.
    complexity=20,
    percent_missing=0.0, # Percentage of missing data (NaN values) in final data set(s).
    causal_graph_config=CausalGraphConfig(
        # Used to complete any unspecified parameters of the CausalGraphConfig.
        graph_complexity=20,
        include_noise=True, # Noise is included in the system. Should always be True.
        max_lag=3, # Maximum possible lag between a parent and child.
        # Minimum possible lag between a parent and child. 0 allows instantaneous effects
        min_lag=1,
        num_targets=1, # Number of target variables.
        num_features=10, # Number of feature variables.
        num_latent=2, # Number of latent variables.
        # Likelihood of an edge. Used to control graph sparsity.
        # Used when prob_<var_type>_parent is undefined for specific <var_type>.
        prob_edge=0.25,
        # Only 1 lagged node of a variable can be a parent of another node.
        max_parents_per_variable=1, # Helps to control graph sparsity.
        # max_<var_type>_parents and max_<var_type>_children help control graph sparsity.
        max_target_parents=2, # Maximum number of parents for a target node.
        max_target_children=0, # Maximum number of children for a target node.
        # Likelihood of an edge between a possible parent and a target variable.
        prob_target_parent=0.2, # Helps to control graph sparsity.
        max_feature_parents=3, # Maximum number of parents for a feature node.
        max_feature_children=2, # Maximum number of children for a feature node.
        max_latent_parents=2, # Maximum number of parents for a latent node.
        max_latent_children=2, # Maximum number of children for a latent node.
        # A latent variable cannot be a direct cause of a target variable.
        allow_latent_direct_target_cause=False,
        # A target variable cannot be a direct cause of another target variable.
        allow_target_direct_target_cause=False,
        # Likelihood of autoregressive relationship for a target variable.
        prob_target_autoregressive=1.0,
        # Likelihood of autoregressive relationship for a feature variable.
        prob_feature_autoregressive=0.8,
        # Likelihood of autoregressive relationship for a latent variable.
        prob_latent_autoregressive=0.5,
        # Likelihood of autoregressive relationship for a noise variable.
        prob_noise_autoregressive=0.1
    ),
    function_config=FunctionConfig(
        # Used to complete any unspecified parameters of the FunctionConfig.
        function_complexity=30,
        # Possible functional dependencies on each edge.
        functions=["linear", "monotonic", "trigonometric"],
        # Likelihood of choosing above functional forms.
        prob_functions=[0.4, 0.5, 0.1]
    ),
    noise_config=NoiseConfig(
        # Used to complete any unspecified parameters of the NoiseConfig.
        noise_complexity=20,
        # Variance of each noise node is drawn from uniform distribution with given range
        noise_variance=[0.01, 0.1],
        # Possible noise distributions.
        distributions=["gaussian", "laplace", "students_t", "uniform"],
        # Likelihood of choosing above distributions.
        prob_distributions=[0.4, 0.2, 0.2, 0.2]
    ),
    runtime_config=RuntimeConfig(
        # Two data sets will be returned. One with 100 and another with 500 observations.
        num_samples=[100, 500],
        # Seeds used to sample from the SCM to produce the two data sets.
        data_generating_seed=[42, 43],
        # Values for latent and noise variables will not be returned.
        return_observed_data_only=True,
        # Data will not be normalized to zero mean and unit variance.
        normalize_data=False
    )
)
```