TRENDDIFF: DECOUPLING INTRINSIC AND MEASURE MENT TRENDS FOR ENHANCED TIME SERIES CAUSAL DISCOVERY

Anonymous authors

Paper under double-blind review

Abstract

Time trends can be classified into intrinsic (real) and measurement (false) trends. There has long been a critical need for techniques to discern them, especially in investment decision-making. In causal discovery, these measurement trends, essentially measurement errors, can significantly impact the performance of algorithms, making it crucial to identify and eliminate them before analysis as well. Recognizing this need, we present a novel algorithm, termed Trend Differentiator (TrendDiff). It is capable of detecting all trend-influenced variables and differentiating between those affected by measurement trends and those displaying intrinsic trends, relying on changing causal module detection and trend-influenced variables' structural properties, respectively. Extensive experiments on synthetic and real-world data demonstrate the efficacy of this approach.

023

006

008 009 010

011

013

014

015

016

017

018

019

021

1 INTRODUCTION

025 026

048

Emerging in the early 1990s, causal discovery algorithms have undergone substantial growth in
 the past decades (Spirtes and Zhang, 2016). These algorithms strive to infer causality from purely
 observational data, serving as valuable instruments in situations where randomized controlled trials
 are rendered impractical due to ethical concerns, financial constraints, and other obstacles. Standing
 at the intersection of explosive data volumes and advancements in computational capabilities, a
 surge in theoretical and applied causal research has ensued. However, the rapid accumulation of data
 presents not only exciting possibilities but significant challenges in causal discovery.

A prevalent challenge is the presence of time trends, frequently encountered in time series. As articulated by Phillips (2005), "No one understands trends, but everyone sees them in data". Prior research has extensively investigated the impact of trends on the efficacy of conventional statistical algorithms (White and Granger, 2011; Wu et al., 2007). Yet their influence on causal discovery remains largely unexplored. Based on the origin, trends can be classified into two distinct categories: intrinsic (real) and measurement (false) trends. In this context, we define the terms "trend", "intrinsic trend", and "measurement trend" as follows:

Definition 1. A trend is a function concerning time within a given data span. Specificly, time trend T = f(t), for $t_{start} \le t \le t_{end}$.

Definition 2. An *intrinsic trend* is inherent to the fundamental mechanisms governing the variables (e.g., global warming, the temperature is really increasing).

Definition 3. A measurement trend is essentially an observation error unique to the recorded values
 (e.g., an observed increase in diagnosed thyroid nodule patients due to enhanced medical techniques,
 despite a stable real incidence rate over time, see Figure 1).

The two types of trends originate from distinct sources, exert disparate impacts, and necessitate differential treatment.

However, there is this impression – time trends, be it an intrinsic trend, or a measurement trend,
 should be removed before analyses – which is not accurate. Undoubtedly, measurement trends, being
 a form of measurement error, necessitate removal. Take constraint-based causal discovery methods,

which rely on conditional independence tests, for example. Figure 2 (a) shows the true causal graph,



Figure 1: The true and observed incidence of the thyroid nodule along time - a typical example of 060 measurement trends.

the variables X_2 and X_3 are not observable, with a measurement trend exhibited in the observed X_2 063 and X_3 . These measurement trends greatly increase the noise in these variables. For X_2 , this rise in 064 noise alters its relationship with its neighbors X_1 and X_3 , weakening the observed dependencies as 065 the measurement trends intensify. Additionally, the inability to accurately observe X_2 's true values 066 hampers its capacity to d-separate X_1 and X_3 , due to the challenge in precisely controlling for X_2 . 067 Analogous phenomena transpire for another measurement-trend variable X_3 . The causal network 068 identified in Figure 2 (b) diverges significantly from the ground truth in such scenarios. To summarize, 069 measurement trends, inherently measurement errors, introduce two issues for constraint-based causal discovery: 1. the dependence between measurement-trend variables and their neighbors weakens with 071 increasing trends; 2. the conditional independence given the measurement-trend variables vanishes, 072 yielding increasing dependence (Scheines and Ramsey, 2016; Zhang et al., 2017). As highlighted 073 in earlier research regarding measurement error in causal discovery, this influence is not limited to constraint-based causal algorithms but also extends to other methodologies, including those based on 074 functional causal models (Zhang et al., 2017). Conversely, intrinsic trends are integral components 075 of the variables and mechanisms, facilitating the identification of underlying causal relationships. 076 Removal of intrinsic trends would decrease the signal-to-noise ratio, leading to lower detection power, 077 and thus should be avoided. Consequently, discerning between intrinsic and measurement trends is 078 crucial before conducting causal discovery analyses. 079

081

061 062

082

083 084

085

880

101

Figure 2: An illustration of how ignoring measurement trends in causal discovery may lead to spurious connections by constraint-based methods. (a) The true causal graph. Variables whose actual values do not match the observed ones are underlined to indicate their true values. Encircled variables signify their unobservability. Here, the circled X_2 and X_3 represent the true, but unobservable, values of the measurement-trend variables. (b) The estimated skeleton based on observed data.

This study introduces the Trend Differentiator (TrendDiff) algorithm, designed to pinpoint variables 090 influenced by trends and differentiate between those affected by measurement trends and those 091 displaying intrinsic trends. It is not only critical in data pre-processing for causal discovery but carries 092 substantial practical importance in decision-making. Discerning true market trends from transient fluctuations is essential for avoiding misallocation of resources in non-viable market opportunities. 094 The ability to accurately identify trend types is key to reducing such investment risks.

- The principal contributions of our work are shown below: 096
- **Problem Formulation**. We parameterize variables with intrinsic and measurement trends using 098 graphical models. While there has already been research regarding measurement errors in causal 099 discovery, no attention has been paid to differentiating time trends. However, as we motivated 100 above, distinguishing intrinsic from measurement trends is of great theoretical and practical value. To the best of our knowledge, this work is the first to formally propose this problem.
- 102 • TrendDiff Algorithm. Employing the method of detecting changing causal modules, we can 103 efficiently identify all variables affected by trends. Subsequently, by harnessing the unique causal structures under intrinsic and measurement trends, we are able to distinguish between them. 104 Integrating these technologies, we present the TrendDiff algorithm, a novel solution specifically 105 designed for the discernment of time trends. 106
- Experimental Validation. We use extensive experimental evaluations, including analyses of a 107 real-world dataset, to demonstrate the robustness and utility of our algorithm.

¹⁰⁸ 2 PARAMETERIZING TRENDS AND RELATED WORK

110 2.1 PARAMETERIZING TIME TRENDS

112 To put intrinsic and measurement trends clearer, we resort to structural equation models (SEMs), where each variable V_i is formulated as a function of its direct causes and an error term ε_i . Here 113 ε_i encapsulates all other unmeasured causes of V_i , with the ε_i values for different variables being 114 mutually independent. Figure 3 (a) depicts a simple causal model, where a direct causal chain is 115 established from variable X_1 , leading to X_2 , and subsequently to X_3 . Each variable is associated 116 with a structural equation, and the model can be parameterized by assigning exact functions to 117 $f(V_i)$, as well as a joint normal distribution to $\varepsilon_1, \varepsilon_2, \varepsilon_3 \sim \mathcal{N}(\mu, \Sigma^2)$. In this case, Σ^2 is diagonal, 118 reflecting the independence among the error terms ε_1 , ε_2 , and ε_3 . Regardless of the functions 119 and free parameter values assigned, the model in Figure 3 (a) exhibits conditional independence: 120 $X_1 \perp \perp X_3 \mid X_2$. In Figure 3 (b), we present the same model as in Figure 3 (a) but with an added 121 intrinsic trend T_2 affecting X_2 . The intrinsic trend T_2 impacts the generation of X_2 and is an inherent 122 part of its underlying mechanisms. In this case, the observed and real values of X_2 are identical. The 123 added intrinsic time trend can go into the causal network through X_2 without altering the original causal structure. Consequently, a trend in X_3 can be observed, which arises due to the influence 124 of T_2 . In Figure 3 (c), we depict the same model but with true values X_2 being "measured" as 125 X_2 , accompanied by a measurement trend T_2 . In this case, the real and observed values of X_2 126 differ. The measurement trend T_2 is present only in the observed X_2 . Due to the collider at X_2 , T_2 127 cannot influence the real values X_2 and is unable to propagate through the original causal network. 128 To summarize, intrinsic and measurement trends are fundamentally the same in form (a function 129 concerning time within a given data span). However, intrinsic trends affect the true value of variables, 130 whereas measurement trends do not. 131

| $X_1 \longrightarrow X_2 \longrightarrow X_3$ | $\begin{array}{c} X_1 \longrightarrow X_2 \longrightarrow X_3 \\ \uparrow \\ (T_2) \end{array}$ | $\begin{array}{c} X_{1} \longrightarrow (\widehat{X}_{2}) \longrightarrow X_{3} \\ \downarrow \\ X_{2} \longleftarrow (\overline{T}_{2}) \end{array}$ |
|---|---|---|
| $X_1 = \varepsilon_1 X_2 = f(X_1, \varepsilon_2)$ | $X_1 = \varepsilon_1 X_2 = f(X_1, T_2, \varepsilon_2)$ | $X_1 = \varepsilon_1$ $\underline{X_2} = f(X_1, \varepsilon_2)$ |
| $X_3 = f(X_2, \varepsilon_3)$ | $X_3 = f(X_2, \varepsilon_3)$ | $X_3 = f(\underline{X_2}, \varepsilon_3)$ $X_2 = f(\underline{X_2}, T_2)$ |
| (a) | (b) | (c) |

Figure 3: Causal models for variables with trends and corresponding equations. (a) A chain graph without trend. (b) X_2 with an intrinsic trend. (c) X_2 with a measurement trend.

2.2 RELATED WORK

140

141 142 143

144

145 Measurement error in causal discovery. Fundamentally, measurement trends represent a problem 146 of measurement error, which adversely affects causal discovery performance. There has already been 147 some research on measurement error in causal discovery. In linear Gaussian contexts, Scheines and 148 Ramsey (2016) parameterized measurement error using SEMs and explored the effect of Gaussian 149 measurement error on fast greedy equivalence search (FGES). Then identifiability conditions in linear 150 Gaussian situations are discussed in Zhang et al. (2017) through factor analysis, with a key identified challenge being the unknown variances of measurement errors E. If known, the covariance matrix 151 of $\hat{\mathbf{X}}$ would be easily accessible and readily used. To address this, Blom et al. (2018) offers an 152 estimate for the upper bound of \mathbf{E} , while Saeed et al. (2020) proposes a consistent partial correlations 153 estimator. In linear non-Gaussian scenarios, Zhang et al. (2018) demonstrates the identifiability of 154 the ordered group decomposition of G, which contains crucial causal information. However, this 155 method depends on over-complete independent component analysis (OICA (Hyvärinen and Oja, 156 2000)), hindered by issues of local optima and high computational complexity (Hoyer et al., 2008; 157 Shimizu et al., 2009), making the practical application of Zhang et al. (2018)'s theoretically sound 158 results challenging. Given this, Dai et al. (2022) defined the Transformed Independent Noise (TIN) 159 condition and exploited it to identify the *ordered group decomposition* by independence tests. 160

161 Particularly regarding the differentiation of intrinsic and measurement trends, this study stands as the first to offer a solution. Distinct from the above-mentioned studies on causal discovery in the presence

3

of measurement error, our research uniquely: 1) distinguish the two types of trends, facilitating data preprocessing to significantly improve data quality. It can be integrated with various analytical tools;
 2) extends its utility beyond merely enhancing causal discovery. It possesses direct and significant practical relevance in investment decision-making.

166 167

168

3 ASSUMPTIONS

169 Instead of relying on the conventional assumption of causal sufficiency, this research adopts a modified 170 concept, termed "pseudo causal sufficiency" (Huang et al., 2020). Traditional causal sufficiency 171 assumes that all common causes (confounders) influencing observed variables are captured in the 172 dataset. However, the occurrence of time trends presents a challenge to this assumption. Time trends 173 typically emerge from intricate, compounded factors, and time trends across various variables might 174 be interlinked owing to certain hidden confounders. These confounders could represent high-level background factors, like economic policies in the stock market. Therefore, rather than assuming the 175 absence of unobserved confounders, our approach operates under the assumption of pseudo causal 176 sufficiency. This assumption signifies that the only unobserved confounders are those inherent in 177 time trends. 178

Assumption 1 (Pseudo causal sufficiency). Any potential confounders can be encapsulated by a mathematically smooth time function. It follows that at each time instance, the values of these confounders are fixed.

Let $\{g_l(C)\}_{l=1}^L$ represent the set of unobserved variables (potentially empty) underlying time trend *T*, in which *C* is assumed to follow a uniform distribution over the considered period. The data points associated with *C* are assumed to be evenly sampled at a specific frequency, making *C* the time index. Furthermore, we define that for each variable V_i , its parents are denoted by PAⁱ, and the local causal processes are represented by the SEM below:

$$V_i = f_i \left(\mathrm{PA}^i, \mathbf{g}^i(C), \theta_i(C), \varepsilon_i \right) \tag{1}$$

Here, $\mathbf{g}^{i}(C) \subseteq \{g_{l}(C)\}_{l=1}^{L}$ signifies the unobserved variables influencing T_{i} (empty when no intrinsic trend is present behind V_{i}), while $\theta_{i}(C)$ represents the effective parameters within the model, also presumed to be functions of C. ε_{i} denotes a disturbance term, independent of C and exhibiting non-zero variance (i.e., the model is non-deterministic). The mutual independence of ε_{i} is also assumed. Note that, the above function (1) is for variables without trends or affected by intrinsic trends only. For those influenced by measurement trends, the real variable and observed variable can be represented by the function (2) and function (3) below, respectively:

$$\underline{V_i} = f_i \left(\mathrm{PA}^i, \varepsilon_i \right) \tag{2}$$

197 198 199

200

196

188

 $V_i = f_i\left(V_i, \mathbf{g}^i(C), \theta_i(C), \varepsilon_i\right) \tag{3}$

In this work, we consider *C* as a random variable, yielding a joint distribution over $\mathbf{V} \cup \{g_l(C)\}_{l=1}^L \cup \{\theta_m(C)\}_{m=1}^n$. We assume that this distribution adheres to the Markov and faithfulness properties with respect to the graph resulting from the subsequent modifications to *G* (*G* represents the causal structure over **V**): add $\{g_l(C)\}_{l=1}^L \cup \{\theta_m(C)\}_{m=1}^n$ to *G*, and for each *i*, add an arrow from each variable in $\mathbf{g}^i(C)$ to V_i and add an arrow from $\theta_i(C)$ to V_i . This extended graph is denoted as G^{aug} . Evidently, *G* is merely the induced subgraph of G^{aug} over **V**.

Assumption 2 (Causal Markov condition and faithfulness). The joint distribution over $\mathbf{V} \cup \{g_l(C)\}_{l=1}^L \cup \{\theta_m(C)\}_{m=1}^n$ is Markov and faithful to the augmented graph G^{aug} .

To enhance clarity and comprehensibility, this work concentrates on instantaneous causal relationships, and the strength of the causal relations does not change over time. Nevertheless, it is worth noting that our framework can be naturally generalized to encompass time-delayed causal relations, akin to how constraint-based causal discovery has been adapted to manage this (see, e.g. (Chu et al., 2008)).

215 We further assume that variables influenced by trends do not function as leaf nodes, where leaf nodes are defined as having no descendants. As depicted in Figure 3, a critical difference exists

between intrinsic and measurement trends in their interactions with the underlying causal network;
intrinsic trends can be incorporated into this network, whereas measurement trends cannot. This
distinction is crucial for our algorithm's capability to differentiate between these trend types. However,
when a trend-influenced variable is a leaf node, its trend, whether intrinsic or measurement, is
unable to integrate into the existing causal network. Therefore, distinguishing between intrinsic and
measurement trends becomes problematic in such cases, as both types exhibit similar characteristics.

223 224 4

225 226

227

228

229 230

231

232 233

4 THE PROPOSED ALGORITHM

In this section, we introduce the proposed algorithm, TrendDiff, designed to identify all variables influenced by trends (phase 1) and distinguish between those affected by measurement trends and those exhibiting intrinsic trends (phase 2).

4.1 Phase 1: Detection of trend-influenced variables and causal structure recovery

In this section, we leverage the detection of changing causal modules to detect variables exhibiting time trends and deduce the causal network for $\mathbf{V} \cup \{C\}$. The core concept hinges on using the (observed) variable C as a surrogate for the unobserved $\{g_l(C)\}_{l=1}^L \cup \{\theta_m(C)\}_{m=1}^n$. In essence, we utilize C to encapsulate the C-specific information. Under the assumptions in Section 3, it is feasible to deploy conditional independence tests on the combined set of $\mathbf{V} \cup \{C\}$ to detect variables with time trends and recover the structure. This is achieved by Algorithm 1 and supported by Theorem 1.

In Algorithm 1, we first construct a complete undirected graph, denoted U_C , which incorporates both 240 C and V. In Step 2 of the algorithm, the decision regarding whether a variable V_i exhibits a time trend 241 is contingent upon the conditional independence between V_i and C, given a subset of other variables. 242 If a time trend is present in V_i , then the module of V_i evolves in conjunction with C. Consequently, 243 the probability distribution of V_i given its non-C parents, namely $P(V_i \mid \text{PA}^i \setminus \{C\})$, will not remain 244 constant across different values of C. As a result, V_i and C are conditionally dependent regardless of 245 any subset of variables. Based on this, we assume that if $V_i \perp L \subset |PA^i \setminus \{C\}$, there should be no 246 time trend in V_i . Conversely, if this assumption does not hold, then we claim to detect variables with 247 time trends. After this, all variables linked to C, referred to as "C-specific variables", are considered 248 to be with time trends. Step 3 aims to discover the skeleton of the causal structure over V. It leverages 249 the results from Step 2: if neither V_i nor V_i is adjacent to C, then C does not need to be involved in 250 the conditioning set. In practice, one may apply any constraint-based search procedures on $\mathbf{V} \cup C$, 251 e.g., SGS (the Spirtes, Glymour and Scheines algorithm) and PC (the Peter-Clark algorithm) (Spirtes et al., 1993). Its (asymptotic) correctness is justified by the following theorem 1. Finally, step 4 is 253 employed to orient the obtained skeleton based on both standard orientation rules and distribution shift. For a comprehensive explanation of the step 4 orientation procedure and the complete proof of 254 Theorem 1, please refer to (Huang et al., 2020). 255

256 257

258 259

260

261

262

263

264

265

Algorithm 1 Detection of Time-trend Variables and Recovery of Causal Structure

- 2: (Detection of time-trend variables) For each *i*, test for the marginal and conditional independence between V_i and *C*. If they are independent given a subset of $\{V_k \mid k \neq i\}$, remove the edge between V_i and *C* in $U_{\mathcal{G}}$.
- 3: (Recovery of causal skeleton) For every i ≠ j, test for the marginal and conditional independence between V_i and V_j. If they are independent given a subset of {V_k | k ≠ i, k ≠ j} ∪ {C}, remove the edge between V_i and V_j in U_G.
- 4: (Orientation) For the obtained skeleton, orient it by standard orientation rules and distribution shift. After the orientation process, we can get the causal network for $\mathbf{V} \cup C$, called G^{phase1} .
- 266 267 268 269

Theorem 1: Given Assumptions made in Section 3, for every $V_i, V_j \in \mathbf{V}$, V_i and V_j are not adjacent in G if and only if they are independent conditional on some subset of $\{V_k \mid k \neq i, k \neq j\} \cup \{C\}$.

296

301

302 303

304

305 306

307

4.2 PHASE 2: UTILIZING STRUCTURAL DIFFERENCES TO DISTINGUISH BETWEEN INTRINSIC AND MEASUREMENT-TREND VARIABLES

After Phase 1, we procured the set of variables exhibiting time trends (those associated with C) as well as the causal network G^{phase1} for $\mathbf{V} \cup C$. In phase 2, we demonstrate that by examining the different structures within causal networks, it is feasible to differentiate variables with intrinsic trends from those influenced by measurement trends.

4.2.1 DISTINGUISH BETWEEN INTRINSIC AND MEASUREMENT TRENDS BY G^{phase1}

As depicted earlier, intrinsic-trend variables do not change the causal network, whereas those variables characterized by measurement trends can induce structural alterations in causal discovery from the observed variables. Next, we delve into how a measurement-trend variable influences the causal structure of *G*^{phase1} and leverage this understanding to partly distinguish between the two trend types.

283 Figure 4 illustrates how a measurement-trend variable alters the output causal structure of Phase 1. In 284 Figure 4 (a), we depict a chain with a measurement trend in X_2 . During Phase 1, the time index C is 285 integrated into our analysis to pinpoint all trend variables. Due to the presence of a measurement trend 286 in X_2 , a connection from C to X_2 is established. Furthermore, based on the conditional independence 287 observed in the actual structure Figure 4(a), we have $T \perp X_3$ and, crucially, $T \not\perp X_3 | X_2$. By 288 extension, because C is a proxy for T, the relationships $C \perp X_3$ and $C \not\perp X_3 | X_2$ should hold. 289 Therefore, X_2 is a collider and the direction is from X_3 to X_2 . The dependency dynamics between 290 X_1 and C follow suit. As a result, the Phase 1 structural outcome for observed variables should be the one shown in **Figure 4** (b). It's worth noting that since the measurement trend T is independent 291 across all variables within the real causal network, no arrow can stem from the measurement-trend 292 variable to other variables in G^{phase1} (cause the observed measurement-trend variable X_2 would 293 always be identified as a collider). In essence, any linkage from a "C-specific variable" to other entities indicates an intrinsic trend. 295



Figure 4: An illustration of how a measurement-trend variable alters the output causal structure of Phase 1. (a) the real structure with a measurement trend in X_2 . (b) the output structure.

In summary, we first employ G^{phase1} to discern intrinsic-trend variables. A "C-specific variable" is deemed to exhibit an intrinsic trend if it possesses any arrow pointing to other variables in G^{phase1} .

4.2.2 DISTINGUISH BETWEEN INTRINSIC AND MEASUREMENT TRENDS BY FURTHER TESTS

308 Having identified certain intrinsic-trend variables based solely on the structure of G^{phase1}, it becomes 309 necessary to undertake additional conditional independence tests for further recognition of other 310 intrinsic-trend variables. As illustrated in Figure 3, the children of time-trend variables serve as 311 critical pivot points in their differentiation process. For variables with intrinsic trends (see Figure 312 3b), there is $T_2 \not\perp X_3$ and $T_2 \perp X_3 | X_2$. Conversely, for variables with measurement trends (see 313 Figure 3c), there is $T_2 \perp \perp X_3$ and $T_2 \not\perp \perp X_3 | X_2$. Thus, the criterion for identifying an intrinsic-trend 314 variable X_2 can be $T_2 \not\perp X_3$ and $T_2 \perp X_3 \mid X_2$. Here T_2 is the trend of X_2 and X_3 is a child of X_2 . Since the trend T_2 is not directly observable in this context. As an alternative, we employ the time 315 index C again, working as a suitable proxy for the unobservable trend. Therefore, the criterion is: 316 $C \not\perp X_3$ and $C \perp X_3 | X_2$. 317

The first row of **Figure 5** illustrates four scenarios of child variables that may arise when screening for the intrinsic-trend variable X_1 . In **Figure 5** (a), no trend is evident in the child variable X_2 , allowing us to easily identify X_1 as an intrinsic-trend variable using our criterion. However, **Figure 5** (b)(c), the child variable X_2 exhibits intrinsic and measurement trends, respectively. Since trends are functions of time, time serves as a confounder (common cause) of trends T_1 and T_2 . In these cases, the path from T_1 to X_2 via the confounder "time" cannot be blocked, as neither "time" nor T_2 is observable (we can obtain a surrogate for T_2 , but it is insufficiently accurate to block the path).



Figure 5: Different scenarios for descendants of intrinsic-trend variables. First row: Four possible cases of intrinsic-trend variable's child nodes in causal networks. (a) Child node without trend. (b) Child node with an intrinsic trend. (c) Child node with a measurement trend. (d) Child node with a trend from other observable nodes. Second row (b-1), (b-2), (c-1), and (c-2): Four possible cases of intrinsic-trend variable's second-order descendant for structure (b) and (c).

Consequently, we cannot distinguish variables with intrinsic trends from those with measurement trends when all child variables have trends. However, if the trend in the child variable X_2 originates from its other observable parent X_3 , as depicted in **Figure 5** (d), the intrinsic-trend variable X_1 is identifiable since we can block the path through "time" by conditioning on X_3 .

347 For structures (b) and (c), first-order descendants (children) do not facilitate distinguishing trend 348 types. However, can second-order descendants provide clarity? Will it help if structures similar to 349 (a) or (d) emerge subsequent to (b) and (c)? The subsequent row illustrates potential second-order 350 descendant structures for both (b) and (c). Although Figure 5 (b-1)(b-2) remain non-identifiable, Figure 5 (c-1)(c-2) can be discerned. The principles behind (c-1) and (c-2) align with those of (a) 351 and (d), namely $C \not\perp X_3$ and $C \perp X_3 | X_1$. It's noteworthy that structures (c-1) and (c-2) essentially 352 represent (a) and (d) but with an added measurement-trend variable subsequent to the intrinsic-trend 353 variable X_1 under examination. Extending this, we can infer that all structures obtained by adding 354 n measurement-trend variables between X_1 and X_2 in structures (a) and (d) can theoretically be 355 identified, where n=0,1,2... 356

In summary, intrinsic-trend variables are discernible in this process only when (1) the intrinsic-trend variable X to be tested possesses at least one descendant variable Y without trends (like structure (a)) or with trends stemming from other observable variables (like structure (d)); and (2) there are no other intrinsic-trend variables on the path from X to Y. Algorithm 2 for Phase 2 is provided in Appendix A.2. By combining Algorithm 1 and 2, we can obtain the proposed TrendDiff algorithm.

362 363

5 EXPERIMENTS

364 365 366

5.1 SIMULATIONS

367 Fixed structure. Synthetic datasets were constructed based on the SEMs described in Appendix 368 Figure 8. Variables X_1, X_2 , and X_7 were specifically designed to show intrinsic trends, while X_3 and 369 X_6 exhibited measurement trends. All trends were modeled as sinusoidal functions with periods w chosen randomly from a uniform distribution Unif([5, 25]). The relationships in the dataset were set 370 371 to be nonlinear, with half of the links following $f^{(1)}(x) = (1 - 4e^{-x^2/2})x$ and the remaining half 372 following $f^{(2)}(x) = (1 - 4x^3 e^{-x^2/2})x$. Noise types (Gaussian, Exponential, Gumbel) and various 373 374 sample sizes (T = 600, 900, 1200, 1500) were incorporated into the simulations. After data generation, 375 we employed TrendDiff to identify variables with intrinsic trends. Additionally, we compared the performance of the PC algorithm on datasets before and after the removal of measurement trends 376 identified by TrendDiff. The effectiveness was quantitatively assessed using F1 score, precision, and 377 recall metrics. Each experimental configuration was repeated across 50 trials to ensure robustness.

378 The results from the fixed structure simulations are illustrated in Figure 6. Figure 6(a) showcases the 379 efficacy of TrendDiff in detecting variables influenced by intrinsic trends. An increase in the length 380 of the dataset correlates with improvements in the algorithm's performance. Specifically, for datasets 381 of 1500 data points or more, the algorithm achieves near-optimal efficiency, with all primary metrics 382 nearing 0.9. Additionally, TrendDiff maintains consistent performance across various types of noise, demonstrating its robustness. Figure 6(b) provides a comparative analysis of the PC algorithm's 383 performance on datasets both before and after the removal of identified measurement trends. The 384 removal of these trends markedly improves the performance of the PC algorithm. 385



Figure 6: Simulation performance. (a) Performance of identifying intrinsic-trend variables. (b) Performance of PC algorithm using data pre and post-elimination of detected measurement trends. 398

Random structure. In addition to the fixed-structure simulations, the TrendDiff algorithm was 400 also assessed using datasets generated from random structures. This random structure simulation 401 maintained the same parameter settings as its fixed structure counterpart, with the exception that the 402 underlying causal network was randomly generated rather than predetermined. The outcomes of 403 the random structure simulations are detailed in the appendix. Specifically, Figure 11 illustrates the 404 algorithm's performance in identifying intrinsic-trend variables across a range of data lengths (T) and 405 noise types. In line with the findings from the fixed structure simulations, TrendDiff demonstrates 406 robustness against various types of noise. Additionally, a consistent trend is observed where the 407 performance of the algorithm improves as the data length increases. Further insights into the stability of our method are illustrated in Figure 12, which demonstrates TrendDiff's resilience in various data 408 dimensions and sparsity levels. Figure 13 evaluates TrendDiff's performance in scenarios involving 409 linear trends. The results show that TrendDiff is particularly proficient in linear-trend scenarios, 410 further highlighting its effectiveness in a broad range of conditions. When tackling practical issues, 411 considering computational complexity becomes essential. The computational efficiency of TrendDiff 412 is demonstrated in Figure 14, which displays the processing times across different data sizes and 413 number of nodes. Notably, the analysis revealed a non-linear increase in processing times with the 414 growth in data length. However, it is important to point out that, even with this escalation for larger 415 datasets, the processing duration stays within a feasible range for practical applications.

417 5.2 REAL DATA

386

387

388 389

390

391

392

394

396

397

399

416

418

We also applied the proposed approach to a real environmental health dataset. This dataset contains 419 daily values of variables regarding air pollution, weather, and sepsis emergency hospital admission in 420 Hong Kong for the period from 2007 to 2018. It is a typical dataset used to assess the interactions 421 between environmental factors and human health. There are pronounced time trends in this data 422 (Figure 7a), rendering it a good application example for the TrendDiff algorithm. In our initial 423 analysis, we applied TrendDiff to determine the intrinsic trend variables within the data. The outcome 424 from Phase 1 indicates that sepsis emergency hospital admissions, CO, O_3 , and NO_2 are variables 425 exhibiting a trend, be it measurement or intrinsic. Subsequently, in phase 2, we differentiated between 426 measurement-trend and intrinsic-trend variables. It was discerned that CO, O_3 , and NO_2 have 427 intrinsic trends while the daily count of sepsis emergency hospital admissions stood out as the sole 428 variable characterized by a measurement trend. This result is consistent with existing evidence. 429 There have been heated discussions in top medical journals about the observed rise in sepsis cases. A prevailing consensus among researchers is that this uptick in sepsis incidences can be largely 430 attributed to the refined definitions and enhanced coding practices for sepsis, rather than the real 431 incidence increase (Rhee et al., 2017; Fleischmann-Struzek et al., 2018). As for the trio of variables

recognized with an intrinsic trend — CO, O_3 , and NO_2 — ample research has been conducted on their time trends. However, none have ascribed these trends to measurement inaccuracies, supporting our results here (Wei et al., 2022).

Beyond simply distinguishing between intrinsic and measurement trends, we also compared causal discovery outcomes before and after the removal of the identified measurement trends. Here we uti-lized the Peter-Clark-momentary-conditional-independence plus (PCMCI+) method, a well-regarded causal discovery algorithm for time series (Runge, 2020). This dataset under scrutiny was a typical environmental health dataset from Hong Kong, with a focus on uncovering environmental factors contributing to sepsis. As depicted in Figure 7(b), the impact of eliminating the identified measure-ment trend is notably significant on the causal analysis results. Our initial analysis, based on the raw data, classified CO and SO_2 as potential mitigating factors against sepsis. However, when the measurement trend associated with sepsis was removed, the analysis showed a different picture. It revealed that temperature alone was a risk factor for sepsis, which is supported by existing evidence (Helbing et al., 2022). Though this analysis did not deal with other factors like seasonality, the notable differences in the findings underscore the critical importance of detecting and correcting measurement trends in causal analysis.



Figure 7: Evaluation of performance using a real-world dataset. (a) Depiction of time series variables. (b) Raw: discovery of structure from raw data by Peter-Clark-momentary-conditional-independence plus (PCMCI+). Detrended: discovery of structure after removal of identified measurement trends by PCMCI+. Here a curved arrow represents a lagged causal relationship, with the lag day shown on the curve. A straight arrow means a contemporaneous association. A straight line terminating in crosses at both ends represents contemporaneous adjacency with unresolved directionality stemming from contradictory orientation rules. The link color refers to the cross-MCI value, which indicates the strength of the relationships. The node color denotes the auto-MCI value, representing how strong the autocorrelation is.

6 CONCLUSION AND DISCUSSIONS

The need to discern intrinsic trends from measurement trends has been a longstanding challenge.
TrendDiff, our innovative algorithm, is tailored to address this difficulty as evidenced by its successful
application in both simulated and real-world scenarios. However, we recognize a few limitations of
this algorithm. Firstly, although we assume trend-influenced variables are non-leaf nodes, differentiating trend types in leaf nodes also holds value. Secondly, in reality, intrinsic and measurement trends
may coexist in variables, a scenario that TrendDiff currently cannot handle. We leave improving
TrendDiff's ability to differentiate trends in leaf nodes and mixed types for future work.

| 486 | REFERENCES |
|-----|--------------|
| 487 | ITEI ERENCES |

| 488 489 | Tineke Blom, Anna Klimovskaia, Sara Magliacane, and Joris M Mooij. An upper bound for random measurement error in causal discovery. <i>arXiv preprint arXiv:1810.07973</i> , 2018. |
|---------------------------------|---|
| 490 491 | Tianjiao Chu, Clark Glymour, and Greg Ridgeway. Search for additive nonlinear time series causal models. <i>Journal of Machine Learning Research</i> , 9(5), 2008. |
| 492 493 494 | Haoyue Dai, Peter Spirtes, and Kun Zhang. Independence testing-based approach to causal discovery under measurement error and linear non-gaussian models. <i>Advances in Neural Information Processing Systems</i> , 35:27524–27536, 2022. |
| 495 496 497 498 | Carolin Fleischmann-Struzek, Antje Mikolajetz, Daniel Schwarzkopf, J Cohen, CS Hartog, M Pletz, P Gastmeier, and K Reinhart. Challenges in assessing the burden of sepsis and understanding the inequalities of sepsis outcomes between national health systems: secular trends in sepsis and infection incidence and mortality in germany. <i>Intensity and mortality</i> 1825, 2018 |
| 500 501 | Dario Lucas Helbing, Leonie Karoline Stabenow, and Reinhard Bauer. Mouse sepsis models: don't forget ambient temperature! <i>Intensive care medicine experimental</i>, 10(1):29, 2022. |
| 502 503 504 505 | Patrik O Hoyer, Shohei Shimizu, Antti J Kerminen, and Markus Palviainen. Estimation of causal effects using linear non-gaussian causal models with hidden variables. <i>International Journal of Approximate Reasoning</i> , 49(2):362–378, 2008. |
| 506 507 508 | Biwei Huang, Kun Zhang, Jiji Zhang, Joseph Ramsey, Ruben Sanchez-Romero, Clark Glymour, and Bernhard Schölkopf. Causal discovery from heterogeneous/nonstationary data. <i>The Journal of Machine Learning Research</i> , 21(1):3482–3534, 2020. |
| 510 511 | Aapo Hyvärinen and Erkki Oja. Independent component analysis: algorithms and applications. <i>Neural networks</i> , 13(4-5):411–430, 2000. |
| 512 513 514 | World Health Organization et al. Global report on the epidemiology and burden of sepsis: current evidence, identifying gaps and future directions. 2020. |
| 515 516 | Peter CB Phillips. Challenges of trending time series econometrics. <i>Mathematics and Computers in Simulation</i> , 68(5-6):401–416, 2005. |
| 517 518 519 520 | Chanu Rhee, Raymund Dantes, Lauren Epstein, David J Murphy, Christopher W Seymour, Theodore J Iwashyna, Sameer S Kadri, Derek C Angus, Robert L Danner, Anthony E Fiore, et al. Incidence and trends of sepsis in us hospitals using clinical vs claims data, 2009-2014. <i>Jama</i> , 318(13): 1241–1249, 2017. |
| 521 522 523 524 525 | Kristina E Rudd, Sarah Charlotte Johnson, Kareha M Agesa, Katya Anne Shackelford, Derrick Tsoi, Daniel Rhodes Kievlan, Danny V Colombara, Kevin S Ikuta, Niranjan Kissoon, Simon Finfer, et al. Global, regional, and national sepsis incidence and mortality, 1990–2017: analysis for the global burden of disease study. <i>The Lancet</i> , 395(10219):200–211, 2020. |
| 526 527 | Jakob Runge. Causal network reconstruction from time series: From theoretical assumptions to practical estimation. <i>Chaos: An Interdisciplinary Journal of Nonlinear Science</i> , 28(7), 2018. |
| 528 529 530 531 | Jakob Runge. Discovering contemporaneous and lagged causal relations in autocorrelated nonlinear time series datasets. In <i>Conference on Uncertainty in Artificial Intelligence</i> , pages 1388–1397. PMLR, 2020. |
| 532 533 534 | Jakob Runge, Peer Nowack, Marlene Kretschmer, Seth Flaxman, and Dino Sejdinovic. Detecting and quantifying causal associations in large nonlinear time series datasets. <i>Science advances</i> , 5 (11):eaau4996, 2019. |
| 535 536 537 | Basil Saeed, Anastasiya Belyaeva, Yuhao Wang, and Caroline Uhler. Anchored causal inference in the presence of measurement error. In <i>Conference on uncertainty in artificial intelligence</i> , pages 619–628. PMLR, 2020. |
| 539 | Richard Scheines and Joseph Ramsey. Measurement error and causal discovery. In <i>CEUR workshop</i> proceedings, volume 1792, page 1. NIH Public Access, 2016. |

- Shohei Shimizu, Patrik O Hoyer, and Aapo Hyvärinen. Estimation of linear non-gaussian acyclic models for latent factors. *Neurocomputing*, 72(7-9):2024–2027, 2009.
- Mervyn Singer, Clifford S Deutschman, Christopher Warren Seymour, Manu Shankar-Hari, Djillali
 Annane, Michael Bauer, Rinaldo Bellomo, Gordon R Bernard, Jean-Daniel Chiche, Craig M
 Coopersmith, et al. The third international consensus definitions for sepsis and septic shock
 (sepsis-3). Jama, 315(8):801–810, 2016.
- Peter Spirtes and Kun Zhang. Causal discovery and inference: concepts and recent methodological advances. In *Applied informatics*, volume 3, pages 1–28. SpringerOpen, 2016.
- Peter Spirtes, Clark Glymour, Richard Scheines, Peter Spirtes, Clark Glymour, and Richard Scheines.
 Discovery algorithms for causally sufficient structures. *Causation, prediction, and search*, pages 103–162, 1993.
- Yaguang Wei, Xinye Qiu, Mahdieh Danesh Yazdi, Alexandra Shtein, Liuhua Shi, Jiabei Yang,
 Adjani A Peralta, Brent A Coull, and Joel D Schwartz. The impact of exposure measurement error
 on the estimated concentration–response relationship between long-term exposure to pm 2.5 and
 mortality. *Environmental Health Perspectives*, 130(7):077006, 2022.
- Halbert White and Clive WJ Granger. Consideration of trends in time series. *Journal of Time Series Econometrics*, 3(1), 2011.
- Zhaohua Wu, Norden E Huang, Steven R Long, and Chung-Kang Peng. On the trend, detrending, and variability of nonlinear and nonstationary time series. *Proceedings of the National Academy of Sciences*, 104(38):14889–14894, 2007.
- Kun Zhang, Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. Kernel-based conditional
 independence test and application in causal discovery. *arXiv preprint arXiv:1202.3775*, 2012.
- Kun Zhang, Mingming Gong, Joseph Ramsey, Kayhan Batmanghelich, Peter Spirtes, and Clark
 Glymour. Causal discovery in the presence of measurement error: Identifiability conditions. *arXiv preprint arXiv:1706.03768*, 2017.
- Kun Zhang, Mingming Gong, Joseph D Ramsey, Kayhan Batmanghelich, Peter Spirtes, and Clark
 Glymour. Causal discovery with linear non-gaussian models under measurement error: Structural
 identifiability results. In *UAI*, pages 1063–1072, 2018.

| 1 | 1 | |
|---|---|--|

594 A APPENDIX / SUPPLEMENTAL MATERIAL

A.1 ASSUMPTIONS

A.1.1 UNDERSTANDING PSEUDO CAUSAL SUFFICIENCY THROUGH A SIMPLIFIED EXAMPLE

In environmental epidemiology, time series analyses are frequently employed to evaluate the immediate effects of air pollution on respiratory health. Here, the exposure of interest is the daily concentration of air pollutants, while the outcome is the daily number of hospital admissions for respiratory diseases. The primary objective is to determine whether day-to-day fluctuations in air pollutant levels influence the day-to-day variation in respiratory-related hospital admissions. Typically, factors like temperature and relative humidity are adjusted for as known confounders in these studies.

Nonetheless, in addition to these controlled variables, there exist other unobserved confounders that
can significantly impact the analysis. These include variables such as socioeconomic status, changes
in policy, and seasonal variations. Given that the duration of this kind of time series studies often
spans several years or more, it's evident that these unobserved confounders could have a substantial
effect on the findings. Therefore, it's crucial to acknowledge and attempt to account for these unseen
factors to ensure the robustness and accuracy of the results.

In this scenario, the notion of pseudo causal sufficiency emerges. Traditionally, it's assumed in such studies that any unobserved confounders are associated with time and can be encapsulated as a smooth function of time, typically represented through splines. By incorporating this spline function into the analysis, the model effectively accounts for unobserved confounders. This assumption has been widely adopted across environmental health studies. However, it's important to note that this example is provided for the sake of understanding pseudo causal sufficiency more easily. In reality, these traditional studies do not explicitly introduce this concept.

618 619

620

596

597 598

A.1.2 NON-TIME-DELAYED RELATIONSHIPS

In real-world scenarios, both non-time-delayed and time-delayed causal relationships are crucial. In 621 theory, all causal relationships involve some form of delay (as we often say, cause precedes effect). 622 However, in real life, due to our limited knowledge of the true nature of various relationships and the 623 limited precision of the data available to us, we often observe relationships that appear to have no 624 delay. For example, the true mechanism by which air pollution affects lung function may involve 625 a delay of one hour—meaning exposure to severe air pollution now might lead to weakened lung 626 function an hour later. However, since we only have access to daily data on air pollution and lung 627 function, this inherently delayed relationship may appear in the data as if there is no delay. Thus, in 628 real life, both non-time-delayed and time-delayed relationships are very common.

In our study, the assumption of no time-delayed relationships is partly made for the sake of clarity and readability of the paper. As this is the initial introduction of intrinsic and measurement trends using graphical models, and the first formal presentation of the differentiation challenge, clarity is paramount. Some of the figures (Figure 5 for example) in the manuscript are already quite complex and difficult to understand without considering delays. If time delays were to be taken into account, the number of nodes would increase manifold, severely affecting the presentation of the problem and the algorithm.

636

637 A.1.3 TREND-INFLUENCED VARIABLES ARE NON-LEAF NODES

638 The impact of measurement trends on causal discovery outcomes is much smaller for leaf nodes 639 compared to non-leaf nodes. This is why we believe that our method remains highly valuable even 640 if it cannot differentiate the trend types in leaf nodes. As described in paragraph 3 and Figure 2, 641 measurement trends introduce two issues for constrained-based causal discovery algorithms: 1. the 642 dependence between measurement-trend variables and their neighbors weakens with increasing 643 trends; 2. the conditional independence given the measurement-trend variable vanishes, yielding 644 increasing dependence. For non-leaf nodes, such as X_2 and X_3 in Figure 2, both these two issues are 645 present. In contrast, leaf nodes like X_1 and X_4 in the same figure are only subject to the first issue. This is because a leaf node, having no children, does not act as a mediator in any relationship. 646

647

| 648 | Δ2 | ALGORITHM 2 |
|-----|------|-------------|
| 649 | 11.2 | ALGORITIM 2 |

| in phase 1, significance threshold α , conditional independence test $Cl(X, Y, Z)$ returning p Ensure: The set of variables exhibiting an intrinsic trend and the set of variables demonstrime measurement trend within V. 1: IntrinsicSet = \emptyset 2: for all $X_i \in "C$ -specific variables" do 3: β = Causal Graph Matrix(G ^{phase1}) 4: $links = \beta_i$. 5: if $1 \in links$ then \triangleright Check if X_i has outgoin 6: Store X_i in IntrinsicSet 7: RestSet = "C-specific variables" - IntrinsicSet 8: for all $X_j \in \text{RestSet do}$ 9: TestNodes = \vee -"C-specific variables" 10: for all $X_j \in \text{TestNodes do}$ 11: JNb = Neighbors(X_j) - X_i 12: for all $n \in \text{Range}(\text{len}(JNb))$ do 13: for all $S_0 \in \text{Combinations}(JNb, n)$ do 14: $S_1 = S_0 + X_i$ 15: $C = \text{Time index}$ 16: $p_0 = Cl(C, X_j, S_0)$ 17: $p_1 = Cl(C, X_j, S_1)$ 18: if $(p_0 < \alpha) \& (p_1 > \alpha)$ then 19: Store X_i in IntrinsicSet 20: MeasurementSet = "C-specific variables" - IntrinsicSet 21: return IntrinsicSet, MeasurementSet 21: return IntrinsicSet, MeasurementSet 21: return IntrinsicSet, MeasurementSet (a) $X_0 = \varepsilon_0$ $X_5 = f(X_3, \varepsilon_5)$ $X_1 = \varepsilon_1$ $X_6 = f(X_5, \varepsilon_6)$ $X_2 = f(X_0, \varepsilon_3)$ $X_8 = f(X_6, \varepsilon_8)$ $X_4 = f(X_0, \varepsilon_3)$ $X_8 = f(X_6, \varepsilon_8)$ $X_4 = f(X_1, X_2, \varepsilon_4)$ $X_9 = f(X_7, \varepsilon_9)$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure All relationships are nonlinear | Require: Dataset V, "C-sp | ecific variables" ide | ntified in pł | nase 1, causal | structure G ^{pha} | se1 identified |
|--|---|--|---|--|---|---|
| Ensure: The set of variables exhibiting an intrinsic trend and the set of variables demonstrimeasurement trend within V. 1: IntrinsicSet = \emptyset 2: for all $X_i \in "C$ -specific variables" do 3: $\beta = \text{Causal Graph Matrix}(G^{\text{phase1}})$ 4: $links = \beta_i$. 5: if $1 \in \text{links then}$ \triangleright Check if X_i has outgoin 6: Store X_i in IntrinsicSet 7: RestSet = "C-specific variables" IntrinsicSet 8: for all $X_i \in \text{RestSet do}$ 9: TestNodes = V - "C-specific variables" 10: for all $X_j \in \text{TestNodes do}$ 11: JNb = Neighbors(X_j) - X_i 12: for all $n \in \text{Range}(\text{len(JNb)})$ do 13: for all $S_0 \in \text{Combinations}(JNb, n)$ do 14: $S_1 = S_0 + X_i$ 15: $C = \text{Time index}$ 16: $p_0 = \text{CI}(C, X_j, S_0)$ 17: $p_1 = \text{CI}(C, X_j, S_1)$ 18: if $(p_0 < \alpha) \& (p_1 > \alpha)$ then 19: Store X_i in IntrinsicSet 20: MeasurementSet = "C-specific variables" - IntrinsicSet 21: return IntrinsicSet, MeasurementSet Note: the Causal Graph Matrix in line 3 outputs a Causal Graph object, where $\beta_{j,i} = 1$ and $\beta_{i,j}$ indicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} = -1$ indicates $i - j; \beta_{i,j} = \beta_{j,i} = 1$ indicates $i \leftrightarrow j$. A.3 FIXED STRUCTURE SIMULATION (a) $X_0 = \varepsilon_0$ $X_5 = f(X_3, \varepsilon_5)$ $X_1 = \varepsilon_1$ $X_6 = f(X_5, \varepsilon_6)$ $X_2 = f(X_0, \varepsilon_3)$ $X_8 = f(X_6, \varepsilon_8)$ $X_4 = f(X_1, X_2, \varepsilon_4)$ $X_9 = f(X_7, \varepsilon_9)$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure. | in phase 1, significance t | hreshold α , condition | nal indepen | dence test CI(. | $X, Y, \boldsymbol{Z})$ retur | ning <i>p</i> -value |
| measurement trend within V. 1: IntrinsicSet = \emptyset 2: for all $X_i \in "C$ -specific variables" do 3: $\beta = \text{Causal Graph Matrix}(\text{Gphase1})$ 4: $links = \beta_i$. 5: if $1 \in links$ then \triangleright Check if X_i has outgoin, 6: Store X_i in IntrinsicSet 7: RestSet = "C-specific variables" - IntrinsicSet 8: for all $X_i \in \text{RestSet do}$ 9: TestNodes = V - "C-specific variables" 10: for all $X_j \in \text{TestNodes do}$ 11: JNb = Neighbors(X_j) - X_i 12: for all $n \in \text{Range}(\text{len}(JNb))$ do 13: for all $S_0 \in \text{Combinations}(JNb, n)$ do 14: $S_1 = S_0 + X_i$ 15: $C = \text{Time index}$ 16: $p_0 = \text{Cl}(C, X_j, S_0)$ 17: $p_1 = \text{Cl}(C, X_j, S_1)$ 18: if $(p_0 < \alpha) \& (p_1 > \alpha)$ then 19: Store X_i in IntrinsicSet 20: MeasurementSet = "C-specific variables" - IntrinsicSet 21: return IntrinsicSet, MeasurementSet Note: the Causal Graph Matrix in line 3 outputs a Causal Graph object, where $\beta_{j,i} = 1$ and $\beta_{i,j}$ indicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} = -1$ indicates $i - j; \beta_{i,j} = \beta_{j,i} = 1$ indicates $i \leftrightarrow j$. A.3 FIXED STRUCTURE SIMULATION (a) $X_0 = \varepsilon_0 \qquad X_5 = f(X_5, \varepsilon_5)$ $X_1 = \varepsilon_1 \qquad X_6 = f(X_5, \varepsilon_6)$ $X_2 = f(X_0, \varepsilon_2) \qquad X_7 = f(X_5, \varepsilon_7)$ $X_3 = f(X_0, \varepsilon_3) \qquad X_8 = f(X_6, \varepsilon_8)$ $X_4 = f(X_1, X_2, \varepsilon_4) \qquad X_5 = f(X_7, \varepsilon_9)$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure. | Ensure: The set of variable | es exhibiting an intri | insic trend a | and the set of | variables dem | nonstrating a |
| 1: IntrinsicSet = \emptyset 2: for all $X_i \in ``C-specific variables'' do 3: \beta = Causal Graph Matrix(G^{phase1})4: links = \beta_i.5: if 1 \in links then \triangleright Check if X_i has outgoin6: Store X_i in IntrinsicSet7: RestSet = `C-specific variables'' - IntrinsicSet8: for all X_i \in RestSet do9: TestNodes = V - `C-specific variables'' 10: for all X_j \in TestNodes do11: JNb = Neighbors(X_j) - X_i12: for all n \in Range(len(JNb)) do13: for all S_0 \in Combinations(JNb, n) do14: S_1 = S_0 + X_i15: C = Time index16: p_0 = Cl(C, X_j, S_0)17: p_1 = Cl(C, X_j, S_1)18: if (p_0 < \alpha) \& (p_1 > \alpha) then19: Store X_i in IntrinsicSet20: MeasurementSet = `C-specific variables'' - IntrinsicSet21: return IntrinsicSet, MeasurementSet11: Note: the Causal Graph Matrix in line 3 outputs a Causal Graph object, where \beta_{j,i} = 1 and \beta_{i,j}11: ndicate i \rightarrow j; \beta_{i,j} = \beta_{j,i} = -1 indicates i - j; \beta_{i,j} = \beta_{j,i} = 1 indicates i \leftrightarrow j.A.3 FIXED STRUCTURE SIMULATION(a) X_0 = \varepsilon_0 \qquad X_5 = f(X_5, \varepsilon_5)X_1 = \varepsilon_1 \qquad X_6 = f(X_5, \varepsilon_6)X_2 = f(X_0, \varepsilon_2) \qquad X_7 = f(X_5, \varepsilon_7)X_3 = f(X_0, \varepsilon_3) \qquad X_8 = f(X_6, \varepsilon_8)X_4 = f(X_1, X_2, \varepsilon_4) \qquad X_9 = f(X_7, \varepsilon_9)Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we notrinsic and measurement trends and generated the simulated data. (b) The visualization structure.$ | measurement trend with | in V. | | | | |
| 2: for all $X_i \in C$ -specific variables do 3: $\beta = Causal Graph Matrix (Gphase1) 4: links = \beta_i.5: fif 1 \in links then \triangleright Check if X_i has outgoin6: Store X_i in IntrinsicSet7: RestSet = "C-specific variables" - IntrinsicSet8: for all X_i \in \text{RestSet do}9: TestNodes = \mathbf{V} \cdot "C-specific variables"10: for all X_j \in \text{TestNodes do}11: JNb = Neighbors(X_j) \cdot X_i12: for all n \in \text{Range}(\text{len}(JNb)) do13: for all S_0 \in \text{Combinations}(JNb, n) do14: S_1 = S_0 + X_i15: C = \text{Time index}16: p_0 = \text{Cl}(C, X_j, S_0)17: p_1 = \text{Cl}(C, X_j, S_1)18: if (p_0 < \alpha) \& (p_1 > \alpha) then19: Store X_i in IntrinsicSet20: MeasurementSet = "C-specific variables" - IntrinsicSet21: return IntrinsicSet, MeasurementSet14: \rightarrow j_i\beta_{i,j} = \beta_{j,i} = -1 indicates i - j_i; \beta_{i,j} = \beta_{j,i} = 1 and \beta_{i,j}17: p_1 = f(X_0, \varepsilon_0)18: X_0 = \varepsilon_0 X_5 = f(X_3, \varepsilon_5)19: X_1 = \varepsilon_1 X_6 = f(X_5, \varepsilon_6)10: X_2 = f(X_0, \varepsilon_2) X_7 = f(X_5, \varepsilon_7)10: X_3 = f(X_0, \varepsilon_3) X_8 = f(X_6, \varepsilon_8)11: X_4 = f(X_1, X_2, \varepsilon_4) X_9 = f(X_7, \varepsilon_9)12: X_4 = f(X_1, X_2, \varepsilon_4) X_9 = f(X_7, \varepsilon_9)13: X_1 = X_1 X_2 = f(X_0, \varepsilon_2) X_7 = f(X_5, \varepsilon_7)14: X_4 = f(X_1, X_2, \varepsilon_4) X_9 = f(X_7, \varepsilon_9)15: X_4 = Hrelatonships are nonlinear$ | 1: IntrinsicSet = \emptyset | | | | | |
| $\beta = \text{Catsat Graph Natrix}(G = x)$ $i \ links = \beta_i.$ $i \ links = \beta_i.$ $\text{Store } X_i \text{ in IntrinsicSet}$ $\text{Store } X_i \text{ castset do}$ $\text{Store } X_i \in \text{RestSet do}$ $\text{Store } X_i \in \text{RestSet do}$ $\text{Store } X_i \in \text{RestSet do}$ $\text{Store all } X_j \in \text{RestSet do}$ $\text{Store all } X_i \in \text{Range}(\text{len(JNb)) \text{ do}$ $\text{Is for all } N \in \text{Range}(\text{len(JNb)) \text{ do}$ $\text{Is for all } N \in \text{Range}(\text{len(JNb)) \text{ do}$ $\text{Is for all } N \in \text{Range}(\text{len(JNb)) \text{ do}$ $\text{Is } for all S_0 \in \text{Combinations}(\text{JNb}, n) \text{ do} \text{Is } for all S_0 \in \text{Combinations}(\text{JNb}, n) \text{ do} \text{Is } y_1 = S_0 + X_i \text{Is } y_1 = Cl(C, X_j, S_0) \text{Is } y_1 = Cl(C, X_j, S_0) \text{Is } y_1 = Cl(C, X_j, S_1) \text{Is } \text{ if } (p_0 < \alpha) \& (p_1 > \alpha) \text{ then} \text{Store } X_i \text{ in IntrinsicSet} \text{ResurementSet} = "C-\text{specific variables" - IntrinsicSet} \text{Note: the Causal Graph Matrix in line 3 outputs a Causal Graph object, where } \beta_{j,i} = 1 \text{ and } \beta_{i,j} \text{indicate } i \rightarrow j; \beta_{i,j} = \beta_{j,i} = -1 \text{ indicates } i - j; \beta_{i,j} = \beta_{j,i} = 1 \text{ indicates } i \leftrightarrow j. \text{A.3 FIXED STRUCTURE SIMULATION \text{(a)} X_0 = \varepsilon_0 \qquad X_5 = f(X_5, \varepsilon_5) \\ X_4 = f(X_1, X_2, \varepsilon_4) \qquad X_9 = f(X_7, \varepsilon_9) \qquad X_0 \qquad X_2 X_4 X_5 X_5 X_4 X_5 X_5 X_5 X_4 X_5 X$ | 2: for all $X_i \in {}^{\circ}U$ -specific | C Variables do | | | | |
| Figure 8: Data structure for fixed-structure simulation. (a) The SEMS Figure 8: Data structure for fixed-structure simulation. (a) The SEMS according to which we distructure All relationships are nonlinear | 5: $\beta = \text{Causar Graph I}$ | viaurix(G ^r | | | | |
| Figure 8: Data structure for fixed-structure simulation. (a) The Set Set Set Set Set Set Set Set Set Se | 5. if $1 \in links$ then | | | ⊳ Chec | k if X: has or | itgoing links |
| 7: RestSet = "C-specific variables" - IntrinsicSet 8: for all $X_i \in \text{RestSet do}$ 9: TestNodes = V - "C-specific variables" 10: for all $X_j \in \text{TestNodes do}$ 11: JNb = Neighbors(X_j) - X_i 12: for all $n \in \text{Range}(\text{len}(JNb))$ do 13: for all $S_0 \in \text{Combinations}(JNb, n)$ do 14: $S_1 = S_0 + X_i$ 15: $C = \text{Time index}$ 16: $p_0 = \text{CI}(C, X_j, S_0)$ 17: $p_1 = \text{CI}(C, X_j, S_1)$ 18: if $(p_0 < \alpha) \& (p_1 > \alpha)$ then 19: Store X_i in IntrinsicSet 20: MeasurementSet = "C-specific variables" - IntrinsicSet 21: return IntrinsicSet, MeasurementSet Note: the Causal Graph Matrix in line 3 outputs a Causal Graph object, where $\beta_{j,i} = 1$ and $\beta_{i,j}$ indicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} = -1$ indicates $i - j; \beta_{i,j} = \beta_{j,i} = 1$ indicates $i \leftrightarrow j$. A.3 FIXED STRUCTURE SIMULATION (a) $X_0 = \varepsilon_0$ $X_5 = f(X_3, \varepsilon_5)$ $X_2 = f(X_0, \varepsilon_2)$ $X_7 = f(X_5, \varepsilon_7)$ $X_3 = f(X_0, \varepsilon_3)$ $X_8 = f(X_6, \varepsilon_8)$ $X_4 = f(X_1, X_2, \varepsilon_4)$ $X_9 = f(X_7, \varepsilon_9)$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure All relationships are nonlinear | 6: Store X_i in Intri | nsicSet | | r chiec | n ii 114 iius ot | ingoing min |
| 8: for all $X_i \in \text{RestSet do}$ 9: TestNodes = V - "C-specific variables" 10: for all $X_j \in \text{TestNodes do}$ 11: JNb = Neighbors $(X_j) - X_i$ 12: for all $n \in \text{Range}(\text{Ien}(JNb))$ do 13: for all $S_0 \in \text{Combinations}(JNb, n)$ do 14: $S_1 = S_0 + X_i$ 15: $C = \text{Time index}$ 16: $p_0 = \text{CI}(C, X_j, S_0)$ 17: $p_1 = \text{CI}(C, X_j, S_1)$ 18: if $(p_0 < \alpha) \& (p_1 > \alpha)$ then 19: Store X_i in IntrinsicSet 20: MeasurementSet = "C-specific variables" - IntrinsicSet 21: return IntrinsicSet, MeasurementSet 21: return IntrinsicSet, MeasurementSet 22: Note: the Causal Graph Matrix in line 3 outputs a Causal Graph object, where $\beta_{j,i} = 1$ and $\beta_{i,j}$ andicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} = -1$ indicates $i - j; \beta_{i,j} = \beta_{j,i} = 1$ indicates $i \leftrightarrow j$. A.3 FIXED STRUCTURE SIMULATION (a) $X_0 = \varepsilon_0$ $X_5 = f(X_5, \varepsilon_5)$ $X_1 = \varepsilon_1$ $X_6 = f(X_5, \varepsilon_7)$ $X_3 = f(X_0, \varepsilon_2)$ $X_7 = f(X_5, \varepsilon_7)$ $X_3 = f(X_0, \varepsilon_3)$ $X_8 = f(X_6, \varepsilon_8)$ $X_4 = f(X_1, X_2, \varepsilon_4)$ $X_9 = f(X_7, \varepsilon_9)$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we ntrinsic and measurement trends and generated the simulated data. (b) The visualization structure. | 7: RestSet = "C-specific v | ariables" - IntrinsicS | Set | | | |
| 9: TestNodes = V - "C-specific variables" 10: for all $X_j \in$ TestNodes do 11: JNb = Neighbors(X_j) - X_i 12: for all $n \in$ Range(len(JNb)) do 13: for all $S_0 \in$ Combinations(JNb, n) do 14: $S_1 = S_0 + X_i$ 15: C = Time index 16: $p_0 = \operatorname{Cl}(C, X_j, S_0)$ 17: $p_1 = \operatorname{Cl}(C, X_j, S_1)$ 18: if $(p_0 < \alpha) \& (p_1 > \alpha)$ then 19: Store X_i in IntrinsicSet 20: MeasurementSet = "C-specific variables" - IntrinsicSet 21: return IntrinsicSet, MeasurementSet 11: Note: the Causal Graph Matrix in line 3 outputs a Causal Graph object, where $\beta_{j,i} = 1$ and $\beta_{i,j}$ 12: indicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} = -1$ indicates $i - j; \beta_{i,j} = \beta_{j,i} = 1$ indicates $i \leftrightarrow j$. A.3 FIXED STRUCTURE SIMULATION (a) $X_0 = \varepsilon_0 \qquad X_5 = f(X_3, \varepsilon_5) \qquad (b)$ $X_1 = \varepsilon_1 \qquad X_6 = f(X_5, \varepsilon_6) \qquad X_7 = f(X_5, \varepsilon_7) \qquad X_3 = f(X_0, \varepsilon_3) \qquad X_8 = f(X_6, \varepsilon_8) \qquad X_4 = f(X_1, X_2, \varepsilon_4) \qquad X_9 = f(X_7, \varepsilon_9)$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure. | 8: for all $X_i \in \text{RestSet do}$ | | | | | |
| 10: for all $X_j \in \text{TestNodes do}$ 11: JNb = Neighbors(X_j) - X_i 12: for all $n \in \text{Range}(\text{len}(JNb))$ do 13: for all $S_0 \in \text{Combinations}(JNb, n)$ do 14: $S_1 = S_0 + X_i$ 15: $C = \text{Time index}$ 16: $p_0 = \text{CI}(C, X_j, S_0)$ 17: $p_1 = \text{CI}(C, X_j, S_1)$ 18: if $(p_0 < \alpha)$ & $(p_1 > \alpha)$ then 19: Store X_i in IntrinsicSet 20: MeasurementSet = "C-specific variables" - IntrinsicSet 21: return IntrinsicSet, MeasurementSet 21: return IntrinsicSet, MeasurementSet 23: fixed STRUCTURE SIMULATION (a) $X_0 = \varepsilon_0$ $X_5 = f(X_3, \varepsilon_5)$ $X_1 = \varepsilon_1$ $X_6 = f(X_5, \varepsilon_6)$ $X_2 = f(X_0, \varepsilon_2)$ $X_7 = f(X_5, \varepsilon_7)$ $X_3 = f(X_0, \varepsilon_3)$ $X_8 = f(X_6, \varepsilon_8)$ $X_4 = f(X_1, X_2, \varepsilon_4)$ $X_9 = f(X_7, \varepsilon_9)$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure. All relationships are nonlinear | 9: TestNodes = \mathbf{V} - " C | -specific variables" | | | | |
| 11: JNb = Neighbors(X_j) - X_i 12: for all $n \in \text{Range}(\text{len}(JNb))$ do 13: for all $S_0 \in \text{Combinations}(JNb, n)$ do 14: $S_1 = S_0 + X_i$ 15: $C = \text{Time index}$ 16: $p_0 = \text{CI}(C, X_j, S_0)$ 17: $p_1 = \text{CI}(C, X_j, S_1)$ 18: if $(p_0 < \alpha) \& (p_1 > \alpha)$ then 19: Store X_i in IntrinsicSet 20: MeasurementSet = "C-specific variables" - IntrinsicSet 21: return IntrinsicSet, MeasurementSet Note: the Causal Graph Matrix in line 3 outputs a Causal Graph object, where $\beta_{j,i} = 1$ and $\beta_{i,j}$ indicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} = -1$ indicates $i - j; \beta_{i,j} = \beta_{j,i} = 1$ indicates $i \leftrightarrow j$. A.3 FIXED STRUCTURE SIMULATION (a) $X_0 = \varepsilon_0 \qquad X_5 = f(X_5, \varepsilon_5)$ $X_1 = \varepsilon_1 \qquad X_6 = f(X_5, \varepsilon_6)$ $X_2 = f(X_0, \varepsilon_2) \qquad X_7 = f(X_5, \varepsilon_7)$ $X_3 = f(X_0, \varepsilon_3) \qquad X_8 = f(X_6, \varepsilon_8)$ $X_4 = f(X_1, X_2, \varepsilon_4) \qquad X_9 = f(X_7, \varepsilon_9)$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure. All relationships are nonlinear | 10: for all $X_j \in \text{TestNo}$ | odes do | | | | |
| 12: for all $n \in \text{Range}(\text{len}(JNb))$ do 13: for all $S_0 \in \text{Combinations}(JNb, n)$ do 14: $S_1 = S_0 + X_i$ 15: $C = \text{Time index}$ 16: $p_0 = \text{CI}(C, X_j, S_0)$ 17: $p_1 = \text{CI}(C, X_j, S_1)$ 18: if $(p_0 < \alpha) \& (p_1 > \alpha)$ then 19: Store X_i in IntrinsicSet 20: MeasurementSet = "C-specific variables" - IntrinsicSet 21: return IntrinsicSet, MeasurementSet Note: the Causal Graph Matrix in line 3 outputs a Causal Graph object, where $\beta_{j,i} = 1$ and $\beta_{i,j}$ indicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} = -1$ indicates $i \rightarrow j; \beta_{i,j} = \beta_{j,i} = 1$ indicates $i \leftrightarrow j$. A.3 FIXED STRUCTURE SIMULATION (a) $X_0 = \varepsilon_0$ $X_5 = f(X_3, \varepsilon_5)$ (b) $X_1 = \varepsilon_1$ $X_6 = f(X_5, \varepsilon_6)$ $X_2 = f(X_0, \varepsilon_2)$ $X_7 = f(X_5, \varepsilon_7)$ $X_3 = f(X_0, \varepsilon_3)$ $X_8 = f(X_6, \varepsilon_8)$ $X_4 = f(X_1, X_2, \varepsilon_4)$ $X_9 = f(X_7, \varepsilon_9)$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure. All relationships are nonlinear | 11: JNb = Neighbor | $\mathbf{s}(X_j)$ - X_i | | | | |
| 13: for all $S_0 \in \text{Combinations(JNb, n)}$ do 14: $S_1 = S_0 + X_i$ 15: $C = \text{Time index}$ 16: $p_0 = \text{CI}(C, X_j, S_0)$ 17: $p_1 = \text{CI}(C, X_j, S_1)$ 18: if $(p_0 < \alpha) \& (p_1 > \alpha)$ then 19: Store X_i in IntrinsicSet 20: MeasurementSet = "C-specific variables" - IntrinsicSet 21: return IntrinsicSet, MeasurementSet Note: the Causal Graph Matrix in line 3 outputs a Causal Graph object, where $\beta_{j,i} = 1$ and $\beta_{i,j}$ indicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} = -1$ indicates $i - j; \beta_{i,j} = \beta_{j,i} = 1$ indicates $i \leftrightarrow j$. A.3 FIXED STRUCTURE SIMULATION (a) $X_0 = \varepsilon_0 \qquad X_5 = f(X_3, \varepsilon_5) \qquad X_7 = f(X_5, \varepsilon_7) \qquad X_3 = f(X_0, \varepsilon_3) \qquad X_8 = f(X_6, \varepsilon_8) \qquad X_4 = f(X_1, X_2, \varepsilon_4) \qquad X_9 = f(X_7, \varepsilon_9)$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure. All relationships are nonlinear | 12: for all $n \in \operatorname{Rang}_{G}$ | ge(len(JNb)) do | \ - | | | |
| Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure. All relationships are nonlinear | 13: for all $S_0 \in C$ | Combinations(JNb, V | n) do | | | |
| Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure. All relationships are nonlinear | 14: $S_1 = S_0$ 15: $C = Tim$ | $+\Lambda_i$ | | | | |
| Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure. All relationships are nonlinear | $\frac{16}{2} \qquad \qquad n_0 = CI($ | $C X \cdot S_{0}$ | | | | |
| Figure 8: Data structure for fixed-structure simulation. (a) The sector $X_i = f(X_1, X_2, \varepsilon_4)$ indicate $X_1 = f(X_1, X_2, \varepsilon_4)$ is a function of the formula formu | 10: $p_0 = CI($ 17: $p_1 = CI($ | (C, X_j, S_0) (C, X_i, S_1) | | | | |
| 19: Store X_i in IntrinsicSet 20: MeasurementSet = "C-specific variables" - IntrinsicSet 21: return IntrinsicSet, MeasurementSet Note: the Causal Graph Matrix in line 3 outputs a Causal Graph object, where $\beta_{j,i} = 1$ and $\beta_{i,j}$ indicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} = -1$ indicates $i - j; \beta_{i,j} = \beta_{j,i} = 1$ indicates $i \leftrightarrow j$. A.3 FIXED STRUCTURE SIMULATION (a) $X_0 = \varepsilon_0 \qquad X_5 = f(X_3, \varepsilon_5) \qquad (b)$ $X_1 = \varepsilon_1 \qquad X_6 = f(X_5, \varepsilon_6) \qquad X_7 = f(X_5, \varepsilon_7) \qquad X_3 = f(X_0, \varepsilon_3) \qquad X_8 = f(X_6, \varepsilon_8) \qquad X_4 = f(X_1, X_2, \varepsilon_4) \qquad X_9 = f(X_7, \varepsilon_9)$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure All relationships are nonlinear | 18: $if(p_0 < $ | α) & $(p_1 > \alpha)$ then | | | | |
| 20: MeasurementSet = "C-specific variables" - IntrinsicSet 21: return IntrinsicSet, MeasurementSet Note: the Causal Graph Matrix in line 3 outputs a Causal Graph object, where $\beta_{j,i} = 1$ and $\beta_{i,j}$ indicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} = -1$ indicates $i - j; \beta_{i,j} = \beta_{j,i} = 1$ indicates $i \leftrightarrow j$. A.3 FIXED STRUCTURE SIMULATION (a) $X_0 = \varepsilon_0$ $X_5 = f(X_3, \varepsilon_5)$ (b) $X_1 = \varepsilon_1$ $X_6 = f(X_5, \varepsilon_6)$ $X_2 = f(X_0, \varepsilon_2)$ $X_7 = f(X_5, \varepsilon_7)$ $X_3 = f(X_0, \varepsilon_3)$ $X_8 = f(X_6, \varepsilon_8)$ $X_4 = f(X_1, X_2, \varepsilon_4)$ $X_9 = f(X_7, \varepsilon_9)$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure. All relationships are nonlinear | <i>z</i> . | | | | | |
| 21: return IntrinsicSet, MeasurementSet Note: the Causal Graph Matrix in line 3 outputs a Causal Graph object, where $\beta_{j,i} = 1$ and $\beta_{i,j}$ indicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} = -1$ indicates $i - j; \beta_{i,j} = \beta_{j,i} = 1$ indicates $i \leftrightarrow j$. A.3 FIXED STRUCTURE SIMULATION (a) $X_0 = \varepsilon_0 \qquad X_5 = f(X_3, \varepsilon_5) \qquad (b)$ $X_1 = \varepsilon_1 \qquad X_6 = f(X_5, \varepsilon_6) \qquad X_7 = f(X_5, \varepsilon_7) \qquad X_3 = f(X_0, \varepsilon_3) \qquad X_8 = f(X_6, \varepsilon_8) \qquad X_4 = f(X_1, X_2, \varepsilon_4) \qquad X_9 = f(X_7, \varepsilon_9)$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure. All relationships are nonlinear | 19: Store | X_i in IntrinsicSet | | | | |
| Note: the Causal Graph Matrix in line 3 outputs a Causal Graph object, where $\beta_{j,i} = 1$ and $\beta_{i,j}$ indicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} = -1$ indicates $i - j; \beta_{i,j} = \beta_{j,i} = 1$ indicates $i \leftrightarrow j$. A.3 FIXED STRUCTURE SIMULATION (a) $X_0 = \varepsilon_0 \qquad X_5 = f(X_3, \varepsilon_5) \qquad (b)$ $X_1 = \varepsilon_1 \qquad X_6 = f(X_5, \varepsilon_6) \qquad X_7 = f(X_5, \varepsilon_7) \qquad X_3 = f(X_0, \varepsilon_3) \qquad X_8 = f(X_6, \varepsilon_8) \qquad X_4 = f(X_1, X_2, \varepsilon_4) \qquad X_9 = f(X_7, \varepsilon_9)$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure. All relationships are nonlinear | 19:Store20:MeasurementSet = " C -s | X_i in IntrinsicSet specific variables" - | IntrinsicSet | | | |
| (a) $X_0 = \varepsilon_0$ $X_5 = f(X_3, \varepsilon_5)$ (b) $X_1 = \varepsilon_1$ $X_6 = f(X_5, \varepsilon_6)$ $X_2 = f(X_0, \varepsilon_2)$ $X_7 = f(X_5, \varepsilon_7)$ $X_3 = f(X_0, \varepsilon_3)$ $X_8 = f(X_6, \varepsilon_8)$ $X_4 = f(X_1, X_2, \varepsilon_4)$ $X_9 = f(X_7, \varepsilon_9)$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we ntrinsic and measurement trends and generated the simulated data. (b) The visualization structure. All relationships are nonlinear | 19: Store 20: MeasurementSet = "C-+ 21: return IntrinsicSet, Me | x_i in IntrinsicSet specific variables" - asurementSet | IntrinsicSet | | | |
| Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure All relationships are nonlinear | 19: Store 20: MeasurementSet = "C- 21: return IntrinsicSet, Me Note: the Causal Graph Matindicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} =$ A.3 FIXED STRUCTURE S | X_i in IntrinsicSet specific variables" - asurementSet rix in line 3 outputs a = -1 indicates i — j SIMULATION | IntrinsicSet Causal Gra ; $\beta_{i,j} = \beta_j$, | ph object, when $a_i = 1$ indicate | ere $oldsymbol{eta}_{j,i}=1$ and \mathbf{s} i $\leftrightarrow j.$ | ad $eta_{i,j}=-1$ |
| $X_{1} = \varepsilon_{1} \qquad X_{6} = f(X_{5}, \varepsilon_{6}) \qquad X_{0} \qquad X_{0} \qquad X_{6} \longrightarrow X_{8}$ $X_{2} = f(X_{0}, \varepsilon_{2}) \qquad X_{7} = f(X_{5}, \varepsilon_{7}) \qquad X_{3} = f(X_{0}, \varepsilon_{3}) \qquad X_{8} = f(X_{6}, \varepsilon_{8}) \qquad X_{4} = f(X_{1}, X_{2}, \varepsilon_{4}) \qquad X_{9} = f(X_{7}, \varepsilon_{9})$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure All relationships are nonlinear | 19: Store 20: MeasurementSet = "C- 21: return IntrinsicSet, Me Note: the Causal Graph Mati indicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} =$ A.3 FIXED STRUCTURE S (a) | x_i in IntrinsicSet specific variables" - asurementSet rix in line 3 outputs a = -1 indicates i — j SIMULATION | IntrinsicSet Causal Gra ; $\beta_{i,j} = \beta_{j,j}$ (b) | ph object, when $i = 1$ indicate | ere $\beta_{j,i} = 1$ and es $i \leftrightarrow j$. $X_{7} = -$ | ad $\beta_{i,j} = -1$ $\longrightarrow X_g$ |
| $X_{2} = f(X_{0}, \varepsilon_{2}) \qquad X_{7} = f(X_{5}, \varepsilon_{7})$ $X_{3} = f(X_{0}, \varepsilon_{3}) \qquad X_{8} = f(X_{6}, \varepsilon_{8})$ $X_{4} = f(X_{1}, X_{2}, \varepsilon_{4}) \qquad X_{9} = f(X_{7}, \varepsilon_{9})$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure. All relationships are nonlinear | 19: Store 20: MeasurementSet = "C- 21: return IntrinsicSet, Me Note: the Causal Graph Mat indicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} =$ A.3 FIXED STRUCTURE S (a) $X_0 = \varepsilon_0$ | x_{i} in IntrinsicSet specific variables" - asurementSet rix in line 3 outputs a = -1 indicates i — j SIMULATION $x_{5} = f(X_{3}, \varepsilon_{5})$ | IntrinsicSet Causal Gra ; $\beta_{i,j} = \beta_j$, (b) | ph object, when $x_i = 1$ indicate | ere $\beta_{j,i} = 1$ and es $i \leftrightarrow j$. $X_{5} \checkmark X_{7} \frown$ | and $\beta_{i,j} = -1$ $\longrightarrow X_g$ |
| $X_{3} = f(X_{0}, \varepsilon_{3}) \qquad X_{8} = f(X_{6}, \varepsilon_{8}) X_{4} = f(X_{1}, X_{2}, \varepsilon_{4}) \qquad X_{9} = f(X_{7}, \varepsilon_{9}) \qquad X_{1} \qquad X_{4}$ Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure All relationships are nonlinear | 19: Store 20: MeasurementSet = "C- 21: return IntrinsicSet, Me Note: the Causal Graph Mati indicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} =$ A.3 FIXED STRUCTURE S (a) $X_0 = \varepsilon_0$ $X_1 = \varepsilon_1$ | $x_{i} \text{ in IntrinsicSet}$ specific variables" - asurementSet rix in line 3 outputs a $= -1 \text{ indicates i } -j$ SIMULATION $X_{5} = f(X_{3}, \varepsilon_{5})$ $X_{6} = f(X_{5}, \varepsilon_{6})$ | IntrinsicSet Causal Gra ; $\beta_{i,j} = \beta_j$, (b) X_0 | ph object, when $x_i = 1$ indicate | ere $\beta_{j,i} = 1$ and es $i \leftrightarrow j$. $X_5 \checkmark X_7 \frown X_6 \frown$ | and $\beta_{i,j} = -1$ $\longrightarrow X_g$ $\longrightarrow X_g$ |
| $X_4 = f(X_1, X_2, \varepsilon_4)$ $X_9 = f(X_7, \varepsilon_9)$ X_1 Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure All relationships are nonlinear | 19: Store 20: MeasurementSet = "C- 21: return IntrinsicSet, Me Note: the Causal Graph Matt indicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} =$ A.3 FIXED STRUCTURE S (a) $X_0 = \varepsilon_0$ $X_1 = \varepsilon_1$ $X_2 = f(X_0, \varepsilon_2)$ | $x_{i} \text{ in IntrinsicSet}$ specific variables" - asurementSet rix in line 3 outputs a $= -1 \text{ indicates i } j$ SIMULATION $X_{5} = f(X_{3}, \varepsilon_{5})$ $X_{6} = f(X_{5}, \varepsilon_{6})$ $X_{7} = f(X_{5}, \varepsilon_{7})$ | IntrinsicSet Causal Gra ; $\beta_{i,j} = \beta_j$, (b) X_0 | ph object, when $x_i = 1$ indicate | ere $\beta_{j,i} = 1$ and so $i \leftrightarrow j$. $X_5 \checkmark X_7 \frown X_6 \frown$ | and $\beta_{i,j} = -1$ $\rightarrow X_g$ $\rightarrow X_g$ |
| Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure All relationships are nonlinear | 19: Store 20: MeasurementSet = "C- 21: return IntrinsicSet, Me Note: the Causal Graph Mattindicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} =$ A.3 FIXED STRUCTURE S (a) $X_0 = \varepsilon_0$ $X_1 = \varepsilon_1$ $X_2 = f(X_0, \varepsilon_2)$ $X_3 = f(X_0, \varepsilon_3)$ | $x_{i} \text{ in IntrinsicSet}$ specific variables" - asurementSet rix in line 3 outputs a $= -1 \text{ indicates i } \text{ j}$ SIMULATION $X_{5} = f(X_{3}, \varepsilon_{5})$ $X_{6} = f(X_{5}, \varepsilon_{6})$ $X_{7} = f(X_{5}, \varepsilon_{7})$ $X_{8} = f(X_{6}, \varepsilon_{8})$ | IntrinsicSet Causal Gra ; $\beta_{i,j} = \beta_j$, (b) X_0 | ph object, when $x_i = 1$ indicates | ere $\beta_{j,i} = 1$ and es $i \leftrightarrow j$. $X_{5} \checkmark X_{7} - X_{6} - X_{6}$ | and $\beta_{i,j} = -1$ $\rightarrow X_g$ $\rightarrow X_g$ |
| Figure 8: Data structure for fixed-structure simulation. (a) The SEMs according to which we intrinsic and measurement trends and generated the simulated data. (b) The visualization structure. All relationships are nonlinear | 19: Store 20: MeasurementSet = "C- 21: return IntrinsicSet, Me Note: the Causal Graph Matindicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} =$ A.3 FIXED STRUCTURE S (a) $X_0 = \varepsilon_0$ $X_1 = \varepsilon_1$ $X_2 = f(X_0, \varepsilon_2)$ $X_3 = f(X_0, \varepsilon_3)$ $X_4 = f(X_1, X_2, \varepsilon_4)$ | $x_{i} \text{ in IntrinsicSet}$ specific variables" - asurementSet trix in line 3 outputs a $= -1 \text{ indicates } i - j$ SIMULATION $X_{5} = f(X_{3}, \varepsilon_{5})$ $X_{6} = f(X_{5}, \varepsilon_{6})$ $X_{7} = f(X_{5}, \varepsilon_{7})$ $X_{8} = f(X_{6}, \varepsilon_{8})$ $X_{9} = f(X_{7}, \varepsilon_{9})$ | IntrinsicSet Causal Gra ; $\beta_{i,j} = \beta_j$, (b) X_0 (| ph object, when $x_i = 1$ indicates | ere $\beta_{j,i} = 1$ and es $i \leftrightarrow j$. $x_5 \checkmark x_7 - x_5 \checkmark x_6 - x_6$ | and $\beta_{i,j} = -1$ $\rightarrow X_g$ $\rightarrow X_g$ |
| intrinsic and measurement trends and generated the simulated data. (b) The visualization structure. All relationships are nonlinear | 19: Store 20: MeasurementSet = "C- 21: return IntrinsicSet, Me Note: the Causal Graph Matindicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} =$ A.3 FIXED STRUCTURE S (a) $X_0 = \varepsilon_0$ $X_1 = \varepsilon_1$ $X_2 = f(X_0, \varepsilon_2)$ $X_3 = f(X_0, \varepsilon_3)$ $X_4 = f(X_1, X_2, \varepsilon_4)$ | $x_{i} \text{ in IntrinsicSet}$ specific variables" - asurementSet rix in line 3 outputs a $= -1 \text{ indicates } i - j$ SIMULATION $X_{5} = f(X_{3}, \varepsilon_{5})$ $X_{6} = f(X_{5}, \varepsilon_{6})$ $X_{7} = f(X_{5}, \varepsilon_{7})$ $X_{8} = f(X_{6}, \varepsilon_{8})$ $X_{9} = f(X_{7}, \varepsilon_{9})$ | IntrinsicSet Causal Gra ; $\beta_{i,j} = \beta_j$, (b) X_0 | ph object, when $x_i = 1$ indicates | ere $\beta_{j,i} = 1$ and es $i \leftrightarrow j$. $x_5 \checkmark x_7 - x_6 - x_6$ | ad $\beta_{i,j} = -1$ $\rightarrow X_g$ $\rightarrow X_8$ |
| structure. All relationships are nonlinear | 19: Store 20: MeasurementSet = "C- 21: return IntrinsicSet, Me Note: the Causal Graph Matindicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} =$ A.3 FIXED STRUCTURE S (a) $X_0 = \varepsilon_0$ $X_1 = \varepsilon_1$ $X_2 = f(X_0, \varepsilon_2)$ $X_3 = f(X_0, \varepsilon_3)$ $X_4 = f(X_1, X_2, \varepsilon_4)$ Figure 8: Data structure for | x_{i} in IntrinsicSet specific variables" - asurementSet rix in line 3 outputs a = -1 indicates i — j SIMULATION $X_{5} = f(X_{3}, \varepsilon_{5})$ $X_{6} = f(X_{5}, \varepsilon_{6})$ $X_{7} = f(X_{5}, \varepsilon_{7})$ $X_{8} = f(X_{6}, \varepsilon_{8})$ $X_{9} = f(X_{7}, \varepsilon_{9})$ fixed-structure simu | IntrinsicSet Causal Gra ; $\beta_{i,j} = \beta_j$, (b) X_0 (| ph object, when $x_i = 1$ indicates $x_1 = 1$ x_2 x_2 x_1 x_2 x_1 The SEMs acc | ere $\beta_{j,i} = 1$ and es $i \leftrightarrow j$. $x_5 \checkmark x_7 - x_6 - x_6$ x_4 ording to which | ad $\beta_{i,j} = -1$ $\rightarrow X_g$ $\rightarrow X_g$ $\rightarrow X_g$ ch we addec |
| succure. The relationships are nonlinear. | 19: Store 20: MeasurementSet = "C- 21: return IntrinsicSet, Me Note: the Causal Graph Mattindicate $i \rightarrow j; \beta_{i,j} = \beta_{j,i} =$ A.3 FIXED STRUCTURE S (a) $X_0 = \varepsilon_0$ $X_1 = \varepsilon_1$ $X_2 = f(X_0, \varepsilon_2)$ $X_3 = f(X_0, \varepsilon_3)$ $X_4 = f(X_1, X_2, \varepsilon_4)$ Figure 8: Data structure for intrinsic and measurement | X_{i} in IntrinsicSet specific variables" - asurementSet rix in line 3 outputs a = -1 indicates i — j SIMULATION $X_{5} = f(X_{3}, \varepsilon_{5})$ $X_{6} = f(X_{5}, \varepsilon_{6})$ $X_{7} = f(X_{5}, \varepsilon_{7})$ $X_{8} = f(X_{6}, \varepsilon_{8})$ $X_{9} = f(X_{7}, \varepsilon_{9})$ fixed-structure simulation trends and generate | IntrinsicSet Causal Gra ; $\beta_{i,j} = \beta_j$, (b) X_0 (b) X_0 (c) d the simul | ph object, when $x_i = 1$ indicates $x_i = 1$ indicates x_2 x_1 The SEMs acclated data. (b | ere $\beta_{j,i} = 1$ and s i $\leftrightarrow j$. $X_{5} \qquad X_{7} - X_{6} - X_{4}$ ording to which) The visualized | and $\beta_{i,j} = -1$ $\rightarrow X_g$ $\rightarrow X_g$ $\rightarrow X_g$ ch we addecontraction of the |

The process to generate simulation data for assessing the TrendDiff algorithm in a fixed structure context involves three primary steps:

 696
 697
 698
 698
 1. Original Structure Acquisition: The baseline fixed structure, void of any trends, is depicted in 697
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 699
 698
 699
 698
 699
 698
 699
 698
 699
 690
 690
 690
 691
 691
 692
 693
 694
 694
 695
 695
 696
 696
 697
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698
 698

699 2. Trend Integration: To incorporate both identifiable structures from Figure 5(a) and (d) into the 700 simulation, intrinsic trends were embedded in variables X_1, X_2 , and X_7 . Additionally, to emulate 701 real-world data characteristics, measurement trends were introduced in X_3 and X_6 . These trends were modeled as smooth functions of time, formulated as $trend = \sin(\frac{w \cdot t}{T})$, where the period w is

722

723

724 725

730 731 732

randomly drawn from a uniform distribution Unif([5, 25]), T represents the data length, and t is the time index.

3. Data Generation and Testing: The final step involved generating simulation data based on the modified structure from the above two steps. All relationships in the data are set to be nonlinear, with 50% of the links using the function $f^{(1)}(x) = (1 - 4e^{-x^2/2})x$ and the other 50% employing $f^{(2)}(x) = (1 - 4x^3e^{-x^2/2})x$. The algorithm's performance was evaluated under a variety of

noise distributions (Gaussian, Exponential, Gumbel) and different sample sizes (T = 600, 900, 1200, 1500). For each scenario, the TrendDiff algorithm was tested using the generated data, with 50 trials conducted in each setting to ensure statistical robustness.

The efficacy of the TrendDiff algorithm is quantitatively assessed using three key metrics: F1 score,
 precision, and recall. These metrics are defined as follows:

Precision: This metric calculates the proportion of true positive outcomes among the total predicted positives. Mathematically, it is expressed as the ratio of the number of true positives (TP) to the sum of true positives and false positives (FP), given by the formula:

$$P = \frac{TP}{TP + FP}$$

Recall: Also known as sensitivity, this metric measures the proportion of actual positives correctly identified. It is calculated as the ratio of true positives to the sum of true positives and false negatives (FN), described by:

$$R = \frac{TP}{TP + FN}$$

F1 Score: The F1 score is the harmonic mean of precision and recall, providing a single metric that balances both. It is particularly useful in situations where there is an uneven class distribution. The F1 score is computed using the formula:

$$F1 = 2 \times \frac{P \times R}{P + R}$$

These metrics offer a comprehensive evaluation of TrendDiff's performance, effectively capturing its accuracy and robustness in various testing scenarios.

735 Our study extends to evaluating the performance of the PC (Peter-Clark) algorithm both before and 736 after the removal of measurement trends identified by TrendDiff. This comparative analysis highlights the effectiveness of our methodology in enhancing causal discovery through data preprocessing. In 737 our approach to handle variables affected by measurement trends, we utilized the Savitzky-Golay 738 filter. This process involves subtracting the trend, as determined by the Savitzky-Golay filter, from 739 the original data to produce detrended data. The Savitzky-Golay filter is a well-known polynomial 740 smoothing technique that fits a polynomial of a specified degree to consecutive data points within a 741 moving window, using linear least squares regression. Once the polynomial is fitted, the filter can 742 provide either a smoothed estimate or the derivative of the fitted function. This filtering method 743 is prevalent in various domains such as analytical chemistry and signal processing, especially for 744 dealing with noisy data. The two primary parameters governing the Savitzky-Golay filter are the 745 window size and the polynomial order. The window size determines the number of data points 746 used for each polynomial fit, thereby influencing the degree of smoothing. On the other hand, the 747 polynomial order specifies the complexity of the model used in the fitting process.

748 Central to our algorithm is the utilization of a nonparametric conditional independence test, which is 749 crucial given the often unknown and highly nonlinear nature of time trends. In this study, we adopted 750 the Kernel-based Conditional Independence test (KCI-test) (Zhang et al., 2012), a method adept 751 at capturing complex nonlinear dependencies. A critical aspect of the KCI-test is the kernel width 752 parameter w, integral to constructing kernel matrices and subsequently influencing the performance 753 of the test. We conducted extensive evaluations to determine the optimal kernel width, varying data lengths T and kernel widths w. These variations in performance, based on different values of T and 754 w, are elucidated in **Figure 9**. The data reveals that as T increases, so does performance efficiency. 755 Notably, a kernel width of w = 0.5 consistently yields impressive results, regardless of the T value. ⁷⁵⁶Our findings are in concordance with the recommendations from the original KCI paper, which ⁷⁵⁷suggest specific kernel width settings based on sample sizes: set w to 0.8 for sample sizes $n \le 200$, ⁷⁵⁸to 0.3 if n > 1200, and to 0.5 in all other instances. In alignment with these guidelines, our study ⁷⁵⁹adopted these kernel width configurations, optimizing our approach for varying data scenarios. This ⁷⁶⁰methodology underscores our commitment to leveraging advanced statistical techniques for accurate ⁷⁶¹and efficient data analysis. The results from the fixed structure simulations with 90% confidence ⁷⁶²interval are illustrated in **Figure 10**.



Figure 9: Parameter choosing results. Performance of our algorithm under different kernel width w with changing data length T.



Figure 10: Simulation performance with 90% confidence interval. (a) Performance of identifying intrinsic-trend variables. (b) Performance of PC algorithm using data pre and post-elimination of detected measurement trends.

A.4 RANDOM STRUCTURE SIMULATION

We also tested our algorithm using simulated data based on random structures. There are three
steps to this process: 1) We generated random graph G from the Erdös-Rényi (ER) random graph
model, with edges added independently with equal probability. The degree, that is, the total number
of edges linked with each node (in + out), is d. Given G, the weights of edges are drawn from

Unif $([-0.6, -0.2] \cup [0.2, 0.6])$ to obtain a weight matrix W_0 . 2) Given W_0 , intrinsic and measurement trends are randomly assigned to variables, with W_0 updated to W. Note that, only intrinsic trend structures like (a) and d in Figure 5 will be generated in this process, which means: a) no trend in leaf nodes; b) variables with trends are not adjacent. 3) Then we sampled $X = W^T X + z \in \mathbb{R}^d$ from noise model. Finally, we generated random datasets $\mathbf{X} \in \mathbb{R}^{n \times d}$ by generating T rows I.I.D. We considered different model setups for noise types, data length T, data dimension, and the degree of sparsity to comprehensively test our algorithm. For each scenario, all metrics precision, recall, and F1 score are computed across all graphs from 50 realizations of the random graph-generating model at data length T in (600, 900, 1200, 1500).

Figure 11 showcases the performance metrics – F1 score, precision, and recall – for identifying intrinsic-trend variables across different data lengths T and noise types. Notably, the method proves robust across noise variations and, consistent with fixed structure results, performance improves with increasing data length. Figure 12 provides further insights into our method's stability, demonstrating its resilience across a range of data dimensions and degrees of sparsity, where dimension is denoted by the number of nodes and sparsity is defined as the degree considering edges in both directions. Figure 13 shows TrendDiff performance on data generated from random structures with linear trends. We measured the identification of intrinsic-trend variables across different data lengths T and noise types in linear-trend scenarios. TrendDiff excels in scenarios with linear trends. Figure 14 provides an analysis of the processing times and peak memory required by TrendDiff for handling different data sizes and number of nodes. The TrendDiff algorithm was executed on a high-performance computing (HPC) system, featuring a single 25-core CPU. A significant finding from this deployment is the non-linear increase in processing times corresponding to the augmentation of data length. Although there is a marked escalation in processing duration for larger datasets, it is essential to emphasize that the timeframes remain within a practical and manageable range for real-world applications. Specifically, for a dataset with 10 variables and a data length of T=1500, the processing time is maintained at approximately 1000 seconds (17min). The peak memory requested is stable.



Figure 11: Performance evaluation on data generated from random structures with varying T and noise type. We measure the identification of intrinsic-trend variables across different data lengths T and noise types using F1 score, precision, and recall. Higher values denote better performance.



Figure 12: Performance evaluation on data generated from random structures with varying sparsity
 and dimension. (a) Performance under different sparsity levels. (b) Performance across varying
 dimensions.



Figure 13: Performance evaluation on data generated from random structures with linear trends. We measure the identification of intrinsic-trend variables across different data lengths T and noise types using F1 score, precision, and recall. Higher values denote better performance.



Figure 14: Run time and peak memory requested by TrendDif with increasing data length T and number of nodes.

A.5 APPLICATION IN REAL-WORLD DATA

875

876

877

878 879

880

883

885

887

888

889

890

891 892 893

894 895

Besides simulation studies, we applied our algorithm to a real-world data set about environmental 896 health as well. The data set contains daily values of variables regarding air pollution, weather, 897 and sepsis emergency hospital admission in Hong Kong. This data set is good for exploring the 898 relationships between environmental factors and sepsis. Sepsis, alternatively referred to as septicemia 899 or blood poisoning, is a life-threatening medical emergency when the dysregulated host response 900 to infection injures its own tissues and organs (Singer et al., 2016). It is one of the leading causes 901 of death and contributes significantly to preventable mortality (Organization et al., 2020). In 2017, 902 11.0 million sepsis-related deaths were reported globally, constituting 20% of all the annual deaths (Rudd et al., 2020). Understanding the relationships between environmental factors and sepsis risk 903 provides a deeper insight into the underlying mechanisms through which environmental factors may 904 predispose, trigger, or exacerbate sepsis conditions. This knowledge is not only pivotal for timely 905 intervention but also offers a foundation for formulating targeted prevention strategies. 906

Data on daily sepsis emergency hospital admissions of Hong Kong were obtained from the Hospital
Authority, which compiles information on all emergency admissions from publicly funded hospitals
that provide 24-hour accident and emergency services and cover 90 percent of hospital beds for
Hong Kong residents. Sepsis cases were identified based on the ninth version of the International
Classification of Diseases (ICD-9: 38), with a total number of 108,831 admissions for a period of
6,543 days spanning from 2007 to 2018.

Hourly concentrations of air pollutants, including carbon monoxide (CO), particulate matter with aerodynamic diameter 2.5m ($PM_2.5$), ozone (O_3), and nitrogen dioxide (NO_2), were obtained from the general air quality monitoring stations in Hong Kong. For daily O_3 concentrations, the maximum 8-hour averages were considered, while 24-hour averages were used for daily concentrations of other air pollutants. The air pollutant data from all monitoring stations were combined to compute city-wide averages for each pollutant. Weather data pertaining to daily average temperature and relative humidity were acquired from the Hong Kong Observatory. The summary statistics of these variables are shown in Table 1.

Table 1: Summary statistics of daily sepsis emergency hospital admissions, air pollution, and weather
 in Hong Kong, 2007-2018^a

| Variables | Mean (SD) | Min | 25th | 50th | 75th | Max |
|----------------------------------|-------------|-------|-------|-------|-------|--------|
| <i>Outcome (daily count)</i> | | | | | | |
| Sepsis | 19(6) | 5 | 15 | 19 | 23 | 39 |
| Air pollution ($\mu g m^{-3}$) | | | | | | |
| CO | 674.9(2322) | 250.0 | 504.0 | 637.3 | 800.6 | 2001.8 |
| PM _{2.5} | 28.9(180) | 4.0 | 14.9 | 24.9 | 38.3 | 138.3 |
| O ₃ | 61.7(366) | 3.3 | 32.6 | 53.6 | 82.7 | 286.5 |
| NO ₂ | 52.3(184) | 4.1 | 39.0 | 49.2 | 62.5 | 162.4 |
| Weather | | | | | | |
| Temperature (°C) | 23.6(52) | 4.9 | 19.3 | 24.8 | 28.2 | 32.4 |
| Humidity (%) | 78.2(105) | 29.0 | 74.0 | 79.0 | 85.0 | 99.0 |

^aAbbreviations: SD = standard deviation; min = minimum value; 25th = 25th percentile; 50th = 50th percentile; 75th = 75th percentile; max = maximum value; CO = carbon monoxide; $PM_{2.5}$ = particulate matter with aerodynamic diameter 2.5m; O_3 = ozone; NO_2 = nitrogen dioxide; Temp. = temperature; Humid. = relative

humidity.

938 939

936

937

940 In the application of our algorithm to real-world datasets, we began by employing the proposed 941 method to systematically identify the sets of variables exhibiting any trends, then focusing specifically 942 on distinguishing between intrinsic-trend variables and measurement-trend variables. These results 943 were rigorously validated against existing research and literature pertaining to trend behaviors and measurement errors in the context of environmental variables and sepsis data. This analysis served 944 to validate the precision of our algorithm. Following this, we applied the "Peter-Clark-momentary-945 conditional-independence plus (PCMCI+)" causal discovery algorithm to the datasets, conducting 946 this procedure both prior to and subsequent to the removal of the identified measurement trends. 947 This two-phase application facilitated a comprehensive comparative analysis, effectively highlighting 948 the impact and advantages of our algorithm in enhancing causal discovery processes. The results 949 from this application demonstrate the utility of our algorithm as a potent data preprocessing tool, 950 significantly aiding in the accuracy and efficacy of subsequent causal analysis. The effectiveness of 951 the algorithm in real-world scenarios, especially in complex fields like environmental studies and 952 medical research, emphasizes its versatility and potential for broader applications.

953 954 Below we detail the PCMCI+ algorithm:

955 PCMCI+ belongs to the so-called constraint-based causal discovery methods family, which is based on conditional independence test(Runge, 2020). Here "PC" refers to the developers Peter and Clark, 956 "MCI" means that the momentary conditional independence (MCI) test idea is added to the traditional 957 PC algorithm, and "+" reminds users that it extends the earlier version of PCMCI to include the 958 discovery of contemporaneous links(Runge et al., 2019). Like other causal graphic models, PCMCI+ 959 works under the general assumptions of the causal Markov condition (each variable in the system 960 is independent of its non-descendants, given its parent variables) and faithfulness (probabilistic 961 information in data emerges not by chance but from causal structures) (Runge, 2018). On top of the 962 general assumptions, two specific assumptions are also requested: causal stationarity (i.e., the causal 963 links hold for all the studied time points) and causal sufficiency (i.e. measured variables include all 964 of the common causes).

PCMCI+ algorithm starts with a skeleton discovery phase, which serves to remove the adjacencies due to indirect paths (mediation) and common causes (confounders). This phase can be divided into lagged stage and contemporaneous stage. The former is to identify lagged potential parents, and the latter is to identify contemporaneous potential parents and optimize identified lagged parents. In the lagged stage, for each variable X_t^j , a superset of lagged ($\tau > 0$) parents $\hat{\beta}_t^-(X_t^j)$ is estimated with the iterative PC1 algorithm. In the contemporaneous stage, we iterate through subsets $S \subset X_t$ of contemporaneous adjacencies and remove adjacencies for all (lagged and contemporaneous) ordered

| 972 973 | pairs $(X_{t-\tau}^i, X_t^j)$ with $X_t^j \in \mathbf{X}_t$ and $X_{t-\tau}^i \in \mathbf{X}_t \cup \widehat{\beta_t}^-(X_t^j)$ if the MCI conditional independence |
|-------------------|---|
| 974 | holds: $(X_{t-\tau}^i \perp X_t^j \mid S, \widehat{\beta_t^-}(X_t^j), \widehat{\beta_{t-\tau}^-}(X_{t-\tau}^i))$. This skeleton discovery phase returns a skeleton |
| 975 | of causal network of undirected relationships among the nodes. |
| 976 977 978 | Next in the orientation phase the contemporaneous links (lagged links can automatically be directed by time order) in the recognized skeleton will be oriented by the collider orientation stage and followed |
| 979 | by the rule orientation stage. In collider orientation process, unshielded triples $X_{t-\tau}^i \to X_t^k \circ - \circ X_t^j$ |
| 980 | (for $\tau > 0$) or $X_t^i \circ - \circ X_t^k \circ - \circ X_t^j$ (for $\tau = 0$) where $X_{t-\tau}^i, X_t^j$ are not adjacent would be oriented |
| 981 | as collider structures if X_t^k is not in the sepset $(X_{t-\tau}^i, X_t^j)$ according to the rule "none". Here sepset |
| 982 | $(\mathbf{v}^i + \mathbf{v}^j)$ more the controlled unichlass that shows that independence of $\mathbf{v}^i + \mathbf{v}^j$ |
| 983 | $(X_{t-\tau}, X_t)$ means the controlled variables when obtaining conditional independence of $X_{t-\tau}, X_t$. |
| 984 985 | Besides the rule "none", another two rules "conservative" and "majority" can also be chosen in this stage. After that, three rules R1, R2, and R3 are followed to orient left links. R1 rule states that |
| 986 | all unambiguous $X_{t-\tau}^i \to X_t^k \circ - \circ X_t^j$ can be oriented as $X_{t-\tau}^i \to X_t^k \to X_t^j$ since there is no |
| 987 | collider left in this stage; in R2 rule, all $X_t^i \to X_t^k \to X_t^j$ structures with $X_t^i \circ - \circ X_t^j$ are oriented |
| 988 | as $X_t^i \to X_t^j$ to avoid circles. Finally, in R3 rule, for all unambiguous $X_t^i \circ - \circ X_t^k \to X_t^j$ and |
| 989 | $X_t^i \circ - \circ X_t^l \to X_t^j$ where X_t^k, X_t^l are independent and $X_t^i \circ - \circ X_t^j$, we orient X_t^i, X_t^j as $X_t^i \to X_t^j$ |
| 990 | to satisfy both the no-collider and no-circle rules. After the orientation process, we leave unoriented |
| 991 | correlations as $\circ - \circ$ and conflicting correlations as $\times - \times$. |
| 992 | For PCMCI+ analysis the Python module "tigramite" (version 5 1 0 3) was used. The main free |
| 993 | parameters of PCMCI+ (in addition to the free parameters of the conditional independence tests) are |
| 994 | the maximum time delay τ_{max} and the significance threshold α_{PC} . We used 3 and 0.05 for these |
| 995 | two parameters, respectively. In the output causal network produced by PCMCI+, a curved arrow |
| 996 | represents a lagged causal relationship, with the lag day shown on the curve. A straight arrow means |

represents a lagged causal relationship, with the lag day shown on the curve. A straight arrow means a contemporaneous association. A conflicting, contemporaneous adjacency "x-x" indicates that the directionality is undecided due to conflicting orientation rules. The link color refers to the cross-MCI value, which indicates the strength of the relationships. The node color denotes the auto-MCI value, representing how strong the autocorrelation is.

1002 A.6 TERMS AND ABBREVIATIONS

1001

1003

1004

Table 2: Glossary of terms

| Term | Definition |
|---------------------------|--|
| Causal discovery | Revealing causal information by analyzing purely observational data under certain assumptions. |
| Time trend | A function concerning time within a given data span. |
| Intrinsic trend | Time trends that are inherent to the fundamental mechanisms governing the variables (real trends). |
| Measurement trend | Time trends that are essentially observation errors unique to the |
| | recorded values (false trends). |
| Causal sufficiency | The absence of unobserved confounders. |
| Pseudo causal sufficiency | Any unmeasured confounding factors influencing the relationship |
| | interested can be adequately represented by a smooth mathemat- |
| | confounders are those inherent in time trends |
| Causal Markov condition | All the relevant probabilistic information that can be obtained |
| Cuusui Markov condition | from the system is contained in its direct causes, or, expressed |
| | differently, if two variables are not connected in the causal graph |
| | given some set of conditions, then they are conditionally indepen- |
| | dent. |
| Causal faithfulness | Independencies in data arise not from coincidence, but rather |
| | from causal structure or, expressed differently, if two variables are |
| | connected by a causal link in the graph. |

| Term | Definition |
|--|---|
| Leaf node Changing causal module | Nodes without any descendants. A component within a causal model or system where the causa relationships can change over time or across different context This concept acknowledges that the dependence between variable are not static and can evolve due to various factors such as shift in underlying mechanisms. |
| | Table 3. Abbreviations |
| Abbreviation | Full description |
| PC algorithm SGS algorithm FGES algorithm PCMCI+ algorithm TrendDiff TIN SEM OICA MCI KCI test PA HPC CO $PM_{2.5}$ NO_2 O_3 Temp. Humid. | The Peter-Clark algorithm The Spirtes-Glymour-Scheines algorithm The Fast Greedy Equivalence Search algorithm The Peter-Clark-momentary-conditional-independence plus algorithm Trend Differentiator Transformed Independent Noise Structural Equation Model Over-complete independent component analysis Momentary conditional independence Kernel-based conditional independence test Parent High-performance computing Carbon monoxide Particulate matter with aerodynamic diameter $\leq 2.5 \mu m$ Nitrogen dioxide Ozone Temperature Relative humidity |