

MAY 12, 2021

Abstract

Despite having a philosophical grounding from empiricism that span some centuries, the algorithmization of causal discovery started only a few decades ago. This formalization of studying causal relationships relies on connections between graphs and probability distributions. In this setting, the task of causal discovery is to recover the graph that best describes the causal structure based on the data we receive. A particular class of causal discovery algorithms, called constraint-based methods rely on Directed Acyclic Graphs (DAGs) as an encoding of Conditional Independence (CI) relations that carry some level of causal information. However, a CI relation such as X and Y being independent conditioned on Z assumes the independence holds for all possible values Z can take, which can tend to be unrealistic in practice where causal relations are often context-specific. In this thesis we aim to develop constraint based algorithms to learn causal structure from Context-Specific Independence (CSI) relations within the discrete setting, where the independence relations are of the form X and Y being independent of $Z = a$ for some a . This is done by using Context-Specific trees, or CStrees for short, which can encode CSI relations.

Sammanfattning

Trots en filosofisk grund från empirism som sträcker sig över några århundraden, började algoritmiseringen av kausal upptäckt bara för några decennier sedan. Denna formalisering av att studera orsakssamband beror på samband mellan grafer och sannolikhetsfördelningar. I den här inställningen är kausal upptäckt att återställa den graf som bäst beskriver kausalstrukturen baserat på de data vi får. En särskild klass av kausala upptäcktsalgoritmer, kallade begränsningsbaserade metoder, är beroende av Directed Acyclic Graphs (DAG) som en kodning av förhållanden med villkorlig självständighet (CI) som bär någon nivå av kausal information. En CI-relation som X och Y är oberoende förutsatt att Z förutsätter att självständigheten gäller för alla möjliga värden Z kan ta, vilket kan vara orealistiskt i praktiken där orsakssamband ofta är sammanhangsspecifikt. I denna avhandling strävar vi efter att utveckla begränsningsbaserade algoritmer för att lära kausalstruktur från Context-Specific Independence (CSI)-relationer inom den diskreta miljön, där självständighetsrelationerna har formen X och Y är oberoende av $Z = a$ för några a . Detta görs med hjälp av sammanhangsspecifika träd, eller kort CStrees, som kan koda CSI-relationer.

Acknowledgements

I would really like to thank Liam for providing me the opportunity to work on a very exciting project, and for being a superb supervisor in every way imaginable. His guidance is without a doubt the secret ingredient that went into this work. A special thanks to George for being on standby for small discussions, and of course to Nikos and Cindy for the virtual study sessions. Also a special thanks to my family, especially my mom, for valuing a good education back when I was younger, and their patience now that I am maybe studying a bit too much. Last but not least, I thank KTH for funding my curiosity.

Contents

1	<i>Introduction</i>	11
1.1	<i>Motivation</i>	11
1.2	<i>Classical methods</i>	11
1.3	<i>Context specific causal discovery</i>	12
1.4	<i>Relevance to machine learning</i>	12
2	<i>Causal Discovery with Directed Acyclic Graphs</i>	15
2.1	<i>The Causal Discovery problem</i>	15
2.2	<i>Direct Acyclic Graphs (DAGs)</i>	16
2.3	<i>Causal Discovery Algorithms for DAGs</i>	19
2.4	<i>Limitations of using DAGs</i>	20
3	<i>Causal Discovery with Context Specific Trees</i>	21
3.1	<i>Context Specific Trees (CStrees)</i>	21
3.2	<i>Learning CStrees from observed data</i>	24
3.3	<i>Understanding high-dimensional CStrees</i>	29
3.4	<i>Model selection for CStrees</i>	38
4	<i>Experiments</i>	41
4.1	<i>Encoding DAGs as CStrees</i>	41
4.2	<i>Recovering the empty context DAG</i>	42
4.3	<i>Synthetic data</i>	43
4.4	<i>Coronary disease data</i>	46
4.5	<i>Mice protein expression data</i>	51
4.6	<i>Supersymmetry data</i>	54
4.7	<i>Vitamin-D dataset</i>	54
4.8	<i>Discussion of empirical performance</i>	57
5	<i>Conclusions</i>	59
5.1	<i>Summary</i>	59
5.2	<i>Future work</i>	59

6	<i>Bibliography</i>	61
7	<i>Appendix A : Implementation remarks</i>	65
8	<i>Appendix B : Extra figures</i>	67

Introduction

1.1 Motivation

At the heart of scientific discovery, and thus of modern society, is our ability to gain knowledge about causal relations based on observations and experimentation. Whilst being an active research area for its own sake, causal discovery has many applications to diverse fields ranging from economics ¹ and genomics ² to the climate sciences ³. In fact, any field involving any form of measurements/observations which are hypothesized to involve causal interactions can benefit from the tools which this research area has to offer, by allowing the user to get some causal model of the system from which their data is generated. Armed with such a causal model, the user can answer more complex queries than they can using just observed data. Such queries form a hierarchy, referred to as Pearl's Causal Hierarchy, which starts from those related to observations, then interventions, then counterfactuals ⁴. In fact it has been recently shown that their strict distinction holds in the measure theoretic sense ⁵.

1.2 Classical methods

Historically speaking the task of causal discovery was achieved by using clever experimental design ⁶. In particular, Randomized Controlled Trials (RCTs) are often called the gold standard in this regard. The common RCT setting involves allocating subjects in the experiment randomly into 2 groups, where one group called the treatment group receives an intervention and the other group called the control group receives no intervention (or treatment), often in the form of a placebo. One of the main limitations of RCTs is that not every system on which we want to infer causal relationships lends to this setting. For example, it is not a good idea to have a RCT to determine

¹ Biwei Huang, Kun Zhang, Mingming Gong, and Clark Glymour. Causal discovery and forecasting in nonstationary environments with state-space models, 2019

² Pengfei Hu, Rong Jiao, Li Jin, and Momiao Xiong. Application of causal inference to genomic analysis: Advances in methodology. *Frontiers in Genetics*, 9:238, 2018

³ Jakob Runge, Sebastian Bathiany, Erik Bollt, Gustau Camps-Valls, Dim Coumou, Ethan Deyle, Clark Glymour, Marlene Kretschmer, Miguel D Mahecha, Jordi Muñoz-Marí, et al. Inferring causation from time series in earth system sciences. *Nature communications*, 10(1):1–13, 2019

⁴ Judea Pearl and Dana Mackenzie. *The Book of Why: The New Science of Cause and Effect*. Basic Books, Inc., USA, 1st edition, 2018

⁵ Duligur Ibeling Thomas Icard, Elias Bareinboim, Juan D. Correa. On pearl's hierarchy and the foundations of causal inference, 2020

⁶ R.A. Fisher. *The design of experiments*. 1935. Oliver and Boyd, Edinburgh, 1935

whether smoking causes lung cancer which would involve forcing the subjects to smoke. This creates a need for more sophisticated methods, namely causal models, to learn causal relationships from both observational and interventional data. In fact one can even use such models to explain the success of RCTs in a formal manner.

1.3 Context specific causal discovery

Much of the causal models in wide use today make use of Directed Acyclic Graphs (DAGs), with the semantics that the directed edges represent causal relationships between the variables, which correspond to the nodes in the graph. As we operate under the assumption that any data we receive is generated by some underlying probability distribution, we relate these graphs to them via conditional independence relations which carry elementary causal information about the system of variables. Many algorithms have been proposed to learn such causal models from observed data in recent decades ⁷, and are in wide use across various domains. They are, however, restrictive in their modelling capacity, since a conditional independence relation by definition is assumed to hold over all possible outcomes of the conditioning set. We call these outcomes contexts, and it is realistic to assume that causal relations depend on context specific information. A number of approaches exist to model such context specific information ⁸, however their representations are not as intuitive as DAGs. Recently, a new representation for context specific causal relations have been proposed, called Context Specific trees (CStrees), which can be seen as a natural generalization of DAGs to the context specific scenario. Namely, CStrees allow for the representation of context specific information as a sequence of DAGs, thus striking a balance between modelling capacity and an intuitive representation. The next natural step is to learn CStrees from data, and this thesis provides the first novel causal discovery algorithms for CStrees to learn from observational data, alongside algorithms to compute the required minimal contexts, which are required for representing CStrees as a sequence of DAGs. This approach is a constraint based approach which generalizes the classic PC algorithm ⁹, alongside applications to both synthetic and real world data.

1.4 Relevance to machine learning

Causal discovery as a subfield of causal modelling contains many ideas which can help in overcoming hard barriers in machine learning. Machine learning can be summarized as the field where practitioners formulate mathematical models of a system of interest,

⁷ David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of machine learning research*, 3(Nov):507–554, 2002; Peter Spirtes and Clark Glymour. An algorithm for fast recovery of sparse causal graphs. *Social Science Computer Review*, 9(1):62–72, 1991; L Solus, Y Wang, and C Uhler. Consistency Guarantees for Greedy Permutation-Based Causal Inference Algorithms. *Biometrika*, 01 2021. asaa104; and Ioannis Tsamardinos, Laura E Brown, and Constantin F Aliferis. The max-min hill-climbing bayesian network structure learning algorithm. *Machine learning*, 65(1):31–78, 2006

⁸ Rodrigo A Collazo, Christiane Görgen, and Jim Q Smith. *Chain event graphs*. CRC Press, 2018; Tomi Silander and Tze-Yun Leong. A dynamic programming algorithm for learning chain event graphs. In *International Conference on Discovery Science*, pages 201–216. Springer, 2013; Peter Thwaites, Jim Q. Smith, and Eva Riccomagno. Causal analysis with chain event graphs. *Artificial Intelligence*, 174(12):889–909, 2010; and Christiane Görgen and Jim Q. Smith. Equivalence classes of staged trees, 2017

⁹ Peter Spirtes and Clark Glymour. An algorithm for fast recovery of sparse causal graphs. *Social Science Computer Review*, 9(1):62–72, 1991

followed by incorporating observed data into this model using various algorithms with the aim of making better predictions about the system. This field has been enjoying significant breakthroughs recently in part due to the availability of a lot of data and faster computers. However, a lot of the work in this field is set in the assumption of independent and identically distributed (i.i.d) data, and ignores information from interventions, domain shifts and temporal structure ¹⁰. As such, there are various problems which still require a causal model, which without it in some cases even give rise to seemingly nuanced paradoxes ¹¹, such as Simpsons paradox ¹², where one might for example have a positive correlation between 2 variables over the whole data, but dividing the samples into further groups would result in a negative correlation within each group. This is not just a theoretical issue, and has been reported in many real life data as well ¹³.

¹⁰ Bernhard Schölkopf. Causality for machine learning, 2019

¹¹ Judea Pearl and Dana Mackenzie. *The Book of Why: The New Science of Cause and Effect*. Basic Books, Inc., USA, 1st edition, 2018

¹² Edward H Simpson. The interpretation of interaction in contingency tables. *Journal of the Royal Statistical Society: Series B (Methodological)*, 13(2):238–241, 1951

¹³ Clifford H. Wagner. Simpson's paradox in real life. *The American Statistician*, 36(1):46–48, 1982

Causal Discovery with Directed Acyclic Graphs

2.1 The Causal Discovery problem

We first provide a formalization of the causal discovery problem. Suppose we have a system of p variables X_1, \dots, X_p which we assume has some underlying probability distribution \mathbb{P} , and from which we have n samples $\{x_1^i, \dots, x_p^i\}_{i=1}^n$. The goal of causal discovery is to recover a structure G that best represents the causal mechanisms of the system. The structure G is often a graph with certain properties that enables it to encode information about the system - this means we must make an assumption that such a structure G exists and it is related to the distribution \mathbb{P} . This information about the system is extracted from the samples we have from the distribution \mathbb{P} - this means we have to make further assumptions to relate information we get from samples in \mathbb{P} to our structure G .

The assumptions to be made are an inevitable artefact of the No Free Lunch theorem ¹ which states that over a uniform distribution over search/learning problems (which includes causal discovery), all algorithms for such problems have equal performance.

There are two common approaches to causal discovery ². The first is constraint-based methods, which treat the problem of finding the structure as a constraint satisfaction problem. One approach to this is to start from a structure where all variables are causally connected then remove connections based on statistical independence from the observed samples. Second is score based methods, which select a causal representation by assigning a score to all possible models, and then choosing a model that minimizes the score. One approach in this direction is to start from a structure where all variables are not causally connected and then proceed to add connections based on how the observed samples give some score, like the Bayesian Infor-

¹ David H. Wolpert. What is important about the no free lunch theorems?, 2020

² Clark Glymour, Kun Zhang, and Peter Spirtes. Review of causal discovery methods based on graphical models. *Frontiers in Genetics*, 10(nil):nil, 2019

mation Criterion (BIC). In this thesis we will mainly be concerned with constraint based methods, particularly in the discrete setting, where we assume the variables in the system can only take discrete values.

2.2 Direct Acyclic Graphs (DAGs)

We now cover some important definitions and concepts related to Directed Acyclic Graphs (DAGs). They are a convenient and informative graphical means of visualizing the direct cause-effect relationships between variables in a system, and the de-facto choice to model causal structures.

Definition 2.1 (DAGs) A Directed Acyclic Graph (DAG) is a directed graph $G = (V, E)$ which has no cycles.

Let $G = (V, E)$ be a graph. We then have the following definitions. A node $u \in V$ is a **parent** of another node $v \in V$ if $(u, v) \in E$, in this case we also say v is a **child** of u . If there is an edge $(u, v) \in E$ or $(v, u) \in E$ we say the nodes u, v are **adjacent**. A set of nodes (u_1, \dots, u_k) , $k \geq 2$ such that $(u_i, u_{i+1}) \in E$ is called a path between u_1 and u_k . This in case, we say u_1 is an **ancestor** of u_k and u_k is a **descendant** of u_1 . If we have a node v that is not a descendant of a node u we say v is a **non-descendant** of u .

For any node $u \in V$ we denote $PA_G(u)$, $CH_G(u)$, $DS_G(u)$, $ND_G(u)$ to be set of parents, children, descendants and non-descendants of u respectively.

Since we are working with discrete probability distributions, we introduce the (open) probability simplex as the space of all possible probability distributions over a set of discrete variables X_1, \dots, X_p whose outcomes are elements of $\mathcal{X} = \prod_{i=1}^p \mathcal{X}_i$.

Definition 2.2 (Probability simplex) Given a finite set \mathcal{X} , The probability simplex on this set is

$$\Delta_{|\mathcal{X}|-1} = \{(f_x : x \in \mathcal{X}) \in \mathbb{R}^{|\mathcal{X}|} : \forall x \in \mathcal{X} f_x > 0, \sum_{x \in \mathcal{X}} f_x = 1\}.$$

Each point in the probability simplex corresponds to a joint distribution over (X_1, \dots, X_p) , and our interest mainly lies to the subset of this space which are connected to structures we can use to model causal relations.

An important concept when relating DAGs to distributions is that of conditional independence, which we define below.

Definition 2.3 (Conditional Independence) Let \mathbb{P} be a distribution with variables X_1, \dots, X_p . Given non-empty subsets $A, B \subset [p]$ and a

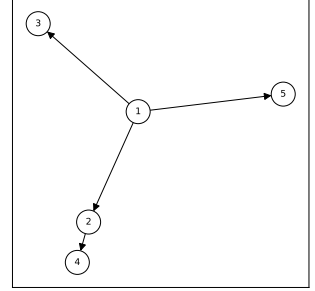


Figure 2.1: Example of a DAG $G = (V, E)$ with $V = \{1, 2, 3, 4, 5\}$ and $E = \{(1, 2), (1, 3), (1, 5), (2, 4)\}$. Here $PA_G(2) = \{1\}$, $DS_G(2) = CH_G(2) = \{1\}$, $ND_G(2) = \{3, 5\}$

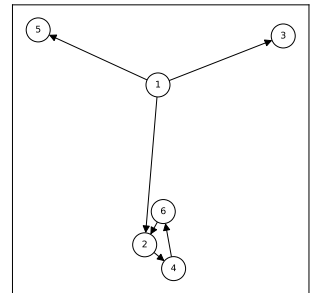


Figure 2.2: This graph is not a DAG since there is a cycle

(possibly empty) subset $S \subset [p]$ such that $\mathbb{P}(X_B, X_S) > 0$ and $A \cap B \cap S = \{\}$, we say the variables X_A and X_B are conditionally independent given S , (denoted $(X_A \perp\!\!\!\perp X_B \mid X_S)$) if $\mathbb{P}(X_A \mid X_B, X_S) = \mathbb{P}(X_A \mid X_S)$ holds for all possible outcomes of X_A, X_B, X_S .

The conditional independence statement $(X_A \perp\!\!\!\perp X_B \mid X_S)$ can be viewed as a ternary relation on X_A, X_B, X_S , and is called a Conditional Independence (CI) relation. This relation formalizes the concept of X_B and X_A not providing any information when we have observed X_S , which is to say, if we already know X_S , knowing X_B does not change the probabilities for X_A , and vice versa.

Using this we can now define the local Markov property which relates distributions to DAGs based on the CI relations encoded by them. As the CI relations have a natural causal interpretation, the local Markov property provides a foundation to relate data generating distributions to DAG representations of a causal system.

[Local Markov property] Let \mathcal{G} be a DAG with nodes $[p]$. A probability distribution \mathbb{P} satisfies the local Markov property with respect to \mathcal{G} if for each node $i \in [p]$, the variable representing that node, X_i is independent of its non-descendants when conditioned on its parents, formally, $(X_i \perp\!\!\!\perp X_{ND_{\mathcal{G}}(i)} \mid X_{PA_{\mathcal{G}}(i)})$

This formalizes the fact that in order to computationally generate data from a DAG \mathcal{G} , the value of each variable $X_i \in \mathcal{X}_i$ depends only on the values of the outcomes of its parents in \mathcal{G} . This means that for a (discrete) distribution \mathbb{P} with p variables satisfying the Local Markov property, the distribution can be encoded with p probability tables which give the probabilities for each X_i taking a value when conditioned on all possible outcomes of its parents. From a storage perspective, this means we have to store $\sum_{i=1}^p |\mathcal{X}_i| \prod_{j \in PA_{\mathcal{G}}(i)} |\mathcal{X}_j|$ which is significantly smaller than having to store all possible probability values which would require one table with $\prod_{i=1}^p |\mathcal{X}_i|$ values. For binary variables assuming d parents for each variable, this is the difference between $p2^{d+1}$ and $p2^p$.

For the purposes of this thesis, it is worth introducing the Ordered Markov property which uses the concept of a linear ordering.³

Definition 2.4 (Ordered Markov Property) Let \mathcal{G} be a DAG and $\pi = \pi_1 \cdots \pi_p$ a causal ordering of \mathcal{G} . A probability distribution \mathbb{P} satisfies the Ordered Markov property with respect to \mathcal{G} if we have $(X_i \perp\!\!\!\perp X_{\{1, \dots, i-1\} \setminus PA_{\mathcal{G}}(i)} \mid X_{PA_{\mathcal{G}}(i)})$

A distribution \mathbb{P} satisfying the local Markov property with respect to a DAG \mathcal{G} is equivalent to that distribution also satisfying the ordered

³ For a DAG \mathcal{G} with p nodes a linear ordering is an ordering of the nodes that respects the directions in \mathcal{G} that is each node i always comes after each $j \in PA_{\mathcal{G}}(i)$. It is also called a topological ordering

Markov property with respect to G and a linear ordering of G .

An important notion in DAGs is that of d-separation and blocked paths.⁴

Definition 2.5 (Blocked path) *Given a DAG G , and a path between nodes $i, j \in \mathbb{V}$, we say the **path is blocked** by a (potentially empty) set of nodes S if either of the following hold:*

- *Along the path there is a triple of nodes (x, s, y) such that $x \rightarrow s \rightarrow y$, $x \leftarrow s \leftarrow y$, or $x \leftarrow s \rightarrow y$ with $s \in S$*
- *Along the path there is a triple of nodes (x, s, y) such that $x \rightarrow s \leftarrow y$ such that $s \notin S$ and no descendants of s are in S .*

Definition 2.6 (d-separation) *Given a DAG G , two (non-empty) sets of nodes X, Y are **d-separated** by a (potentially empty) set of nodes S in G , denoted $(X \perp\!\!\!\perp_G Y \mid S)$ if all paths between every node in X and every node in Y are blocked by S .*

The notion of d-separation relates DAGs to probability distributions from the following theorem.

Definition 2.7 (Global Markov property) *Given a distribution \mathbb{P} that satisfies the local Markov property with a DAG G , we have that for any (non-empty) sets A, B and (possibly empty) set S , $(X_A \perp\!\!\!\perp_G X_B \mid X_S) \implies (X_A \perp\!\!\!\perp_{\mathbb{P}} X_B \mid X_S)$*

An important result is the following that the above notions are indeed equivalent⁵.

Theorem 2.8 (Markov theorems for DAGs) *Given a distribution \mathbb{P} over X_1, \dots, X_p and a DAG G over p nodes, the following are equivalent*

- \mathbb{P} is Markov to G i.e. $\mathbb{P}(X_1, \dots, X_p) = \prod_{i=1}^p \mathbb{P}(X_i \mid X_{PA_G(i)})$
- \mathbb{P}, G satisfy the local Markov property
- \mathbb{P}, G satisfy the ordered Markov property
- \mathbb{P}, G satisfy the global Markov property

If \mathbb{P} satisfies the local Markov property with respect to G and has a probability density with respect to a product measure, we say \mathbb{P} is Markov with respect to G , or equivalently, G is an Independence map (I-MAP) of \mathbb{P} ⁶.

Thus DAGs can be used to store Conditional Independence (CI) relations. More importantly, d-separation encodes the complete set of CI relations satisfied by all distributions Markov to a DAG, i.e.

⁴ A path between 2 nodes is any set of edges connecting them irrespective of the direction.

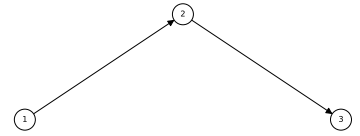


Figure 2.3: Chain

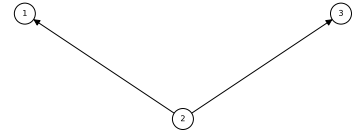


Figure 2.4: Fork/Common cause

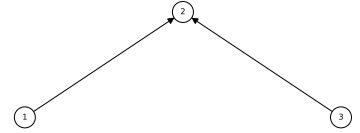


Figure 2.5: V-structure/ Collider/Immortality

⁵ Eliana Duarte and Liam Solus. Algebraic geometry of discrete intervention models, 2020

⁶ Steffen L Lauritzen. *Graphical models*, volume 17. Clarendon Press, 1996

distributions that are Markov to a DAG G and satisfy **exactly** the CI relations encoded by d-separation exist ⁷.

It is also possible to have 2 DAGs that encode the same CI relations, in which case we say that they are both in the same Markov Equivalence Class (MEC), and we say they are Markov Equivalent. MECs can be characterized by the following theorem ⁸.

Theorem 2.9 (Characterization of MECs) *Two DAGs G_1 and G_2 are Markov Equivalent if and only if they have the same skeleton (underlying undirected edges) and v -structures, where a v -structure is a triple of nodes (i, j, k) with edges $i \rightarrow j \leftarrow k$ and i, k do not share an edge.*

For example, the Chain and Fork graphs from the previous page belong to the same Markov Equivalence class.

2.3 Causal Discovery Algorithms for DAGs

Theorem 2.8 suggests that we can make use of CI testing on a distribution \mathbb{P} to learn a DAG G . However, the distribution \mathbb{P} may contain CI relations not encoded in the DAG, thus we make the following assumption.

Definition 2.10 (Faithfulness) *A probability distribution \mathbb{P} is faithful to a DAG G if it entails only the CI relations encoded by the d-separations in the DAG.*

Under the faithfulness assumption, the global Markov property holds both ways. It should be noted that faithful distributions exist ⁹, and the set of distributions that are not faithful to a dag G have measure 0 ¹⁰, which suggests that in theory this is not a very restrictive assumption.

One of the first practical algorithms which make use of the theory above is the PC algorithm, ¹¹ which is a constraint based causal discovery algorithm that relies of the characterization of DAGs in Theorem 2.9 and the faithfulness assumption to find a DAG in the MEC of the true causal DAG. The algorithm starts from a complete graph and runs conditional independence tests to first find the DAG skeleton and then proceeds to direct the edges whenever possible. The output of the PC algorithm is a Completed Partially Directed Acyclic Graph (CPDAG) ¹², which acts as a representation for the Markov Equivalence class. A Partially Directed Acyclic Graph (PDAG) is a graph where some edges are directed and some are undirected and there is no cycle in the direction of the directed edges and any direction of the undirected edges. A PDAG a is Complete PDAG (CPDAG) if every directed edge exists also in every DAG in the

⁷ Christopher Meek. Strong completeness and faithfulness in bayesian networks. *arXiv preprint arXiv:1302.4973*, 2013; and

⁸ Tom S. Verma and Judea Pearl. On the equivalence of causal models. *CoRR*, abs/1304.1108, 2013

⁹ Christopher Meek. Strong completeness and faithfulness in bayesian networks. *arXiv preprint arXiv:1302.4973*, 2013

¹⁰ Caroline Uhler, Garvesh Raskutti, Peter Bühlmann, and Bin Yu. Geometry of the Faithfulness Assumption in Causal inference. *The Annals of Statistics*, 41(2):436 – 463, 2013

¹¹ P. Spirtes, C. Glymour, and R. Scheines. *Causation, Prediction, and Search*. MIT press, 2nd edition, 2000; and Markus Kalisch and Peter Bühlmann. Estimating high-dimensional directed acyclic graphs with the pc-algorithm. *Journal of Machine Learning Research*, 8(22):613–636, 2007

¹² Christopher Meek. Causal inference and causal explanation with background knowledge, 2013

Markov Equivalence class of the DAG and for every undirected edge between nodes i, j there exists a DAG with the edge $i \rightarrow j$ and a DAG with $j \rightarrow i$ in the equivalence class. CPDAGs are also sometimes called essential graphs¹³.

2.4 Limitations of using DAGs

DAGs are a simple and informative structure for causal discovery, however their ability to only encode CI relations is a limitation. This is because the CI relation $(X_A \perp\!\!\!\perp_{\mathbb{P}} X_B \mid X_S)$ implies that X_A and X_B are independent for all possible outcomes of X_S , which in some cases might be too strong of an assumption. A generalization of such relations is Context Specific Independence (CSI) relations, defined below.

Definition 2.11 (Context Specific Independence) *Let \mathbb{P} be a distribution with variables X_1, \dots, X_p with a state space $\mathcal{X} = \prod_{i=1}^p \mathcal{X}_i$. Given non-empty subsets $A, B \subset [p]$ and (possibly empty) subsets $S, C \subset [p]$ and $x_C \in \prod_{i \in C} \mathcal{X}_i$ such that $\mathbb{P}(X_B, X_S, X_C = x_C) > 0$ and $A \cap B \cap S \cap C = \{\}$, we say the variables X_A and X_B are conditional independent given S , in the context $X_C = x_C$ (denoted $(X_A \perp\!\!\!\perp_{\mathbb{P}} X_B \mid X_S)$) if $\mathbb{P}(X_A \mid X_B, X_S, X_C = x_C) = \mathbb{P}(X_A \mid X_S, X_C = x_C)$ holds for all possible outcomes of X_A, X_B, X_S .*

In the next chapter we introduce Context Specific Trees (CStrees) which can encode such relations, and thus provide a structure that can capture the context specific information glossed over in DAGs.

¹³ Steen A. Andersson, David Madigan, and Michael D. Perlman. A characterization of markov equivalence classes for acyclic digraphs. *The Annals of Statistics*, 25(2), April 1997

Causal Discovery with Context Specific Trees

One intuition is that to capture context specific relations one needs to make use of a structure that explicitly represents separate outcomes of a distribution. Typically in high school some might have encountered the use of trees to model small probabilistic systems, and they fully include all possible outcomes involved, and serve as an important tool to compute probabilities for relevant events. As we will see in this chapter, this is a good way to approach the problem of encoding context information as well.

3.1 Context Specific Trees (CStrees)

Before defining CStrees we start by defining staged trees, which contain CStrees as a subset. Both of these are rooted trees.¹

Definition 3.1 (Staged trees) Let $\mathbb{T} = (\mathbb{V}, \mathbb{E})$ be a rooted tree, \mathbb{L} a finite set of labels for the edges, and $\theta : \mathbb{E} \rightarrow \mathbb{L}$ a labelling of the edges. Let $E_{\mathbb{T}}(v) = \{v \rightarrow w \in \mathbb{E} : w \in CH_{\mathbb{T}}(v)\}$, i.e. the set of edges coming out of v in \mathbb{T} . The pair (\mathbb{T}, θ) is a staged tree if

- $\forall v \in \mathbb{V}$ we have $|\theta(E_{\mathbb{T}}(v))| = |E_{\mathbb{T}}(v)|$
- $\forall v, w \in \mathbb{V}$ we have that both $\theta(E_{\mathbb{T}}(v))$ and $\theta(E_{\mathbb{T}}(w))$ are either equal or disjoint

This can be thought of as a probability tree where each edge represents a probability value, and the probabilities coming out of all edges from any given node sum to 1. More formally, first define the space of canonical parameters of the staged tree (\mathbb{T}, θ) as

$$\Theta_{\mathbb{T}} = \{x \in \mathbb{R}^{|\mathbb{L}|} : \forall e \in \mathbb{E}, x_{\theta(e)} \in (0, 1), \forall v \in \mathbb{V}, \sum_{e \in E_{\mathbb{T}}(v)} x_{\theta(e)} = 1\}.$$

Given the probability simplex $\Delta_{|\mathcal{X}|-1}$ and letting $\mathbf{i}_{\mathbb{T}}$ be the set of all leaves of the staged tree \mathbb{T} the staged tree model is defined as below.

¹ A rooted tree $\mathbb{T} = (\mathbb{V}, \mathbb{E})$ is a directed graph whose skeleton is a tree and there exists a unique node r such that $PA_{\mathbb{T}}(r) = \{\}$ which is called the root.

Definition 3.2 (Staged tree models) The staged tree model $\mathbb{M}_{(\mathbb{T}, \theta)}$ is the image of the map $\varphi_{\mathbb{T}} : x \rightarrow f_v := \left(\prod_{e \in E_{\mathbb{T}}(\lambda(v))} x_{\theta(e)} \right)_{v \in \mathbf{i}_{\mathbb{T}}}$

Thus given variables X_1, \dots, X_p , a causal ordering π , the staged tree for this with levels ² $L_1, \dots, L_p \sim X_{\pi_1}, \dots, X_{\pi_p}$, each path from the root to the leaf defines a sequence of events $x_1, x_1 x_2, \dots, x_1 \cdots x_p$ where $x_i \in \mathcal{X}_{\pi_i}$. Since for the edge $e = ((x_1 \cdots x_k), (x_1 \cdots x_k x_{k+1}))$ we have $x_{\theta(e)} = \mathbb{P}(x_{k+1} | x_1 \cdots x_k)$, the product in Definition 3.2 does indeed result in $\mathbb{P}(v_1, \dots, v_p)$ for each $v \in \mathbf{i}_{\mathbb{T}}$ by the chain rule in probability. ³

The important characteristic of staged trees are the stages.

Definition 3.3 (Stages) Given a staged tree (\mathbb{T}, θ) , we say two nodes v, w are in the same stage if and only if $\theta(E_{\mathbb{T}}(v)) = \theta(E_{\mathbb{T}}(w))$

Stages are represented by colours, and when a stage contains a single node, it is coloured white. Staged tree models generalize DAG models, i.e. distributions represented by DAGs, however they are perhaps too general, in the sense that despite allowing for the representation of context specific information, they do not admit a intuitive representation of the causal structure. This creates the need for a structure that generalizes DAG models **and** admits an intuitive representation. The recently proposed subclass of staged trees, known as CStrees allow for this.

Definition 3.4 (CStrees) Let \mathcal{X}_i denote the state space of some variable X_i with $\mathcal{X} = \prod_{i=1}^p \mathcal{X}_i$, and (\mathbb{T}, θ) be a staged tree with levels L_1, \dots, L_p corresponding to variables $X_{\pi_1}, \dots, X_{\pi_p}$ where $\pi = \pi_1 \dots \pi_p$ is the causal ordering of the variables. A CStree is a staged tree (\mathbb{T}, θ) where each level of the tree corresponds to some variable and such that

- It is compatibly labelled, i.e. $\forall x_{\pi_k} \in \mathcal{X}_{\pi_k}$ we have $\theta(x_{\pi_1} \dots x_{\pi_{k-1}} \rightarrow x_{\pi_{k-1}} x_{\pi_k}) = \theta(y_{\pi_1} \dots y_{\pi_{k-1}} \rightarrow y_{\pi_{k-1}} x_{\pi_k})$ whenever $x_{\pi_1} \dots x_{\pi_{k-1}}$ and $y_{\pi_1} \dots y_{\pi_{k-1}}$ are in the same stage
- (**CStree property**) Each stage $S_i \subset L_k$ of the tree has a fixed context, i.e. $\exists C_i \subset [k]$ and the fixed outcome $x_{C_i} \in \mathcal{X}_{C_i}$, where the stages are the union over the variables beside those in C_i , i.e. if $Y_i = [k] \setminus C_i$ then $S_i = \bigcup_{x_{Y_i} \in \mathcal{X}_{Y_i}} \{x_{C_i} x_{Y_i}\}$

Given a CStree \mathbb{T} and a causal ordering π , each node in level L_k corresponds to an outcome of the sequence of variables $X_{\pi_1}, \dots, X_{\pi_k}$. Each edge coming into each node in L_k is of the form $(x_1 \cdots x_{k-1}, x_1 \cdots x_k)$ represents $P(x_k | x_1 \cdots x_{k-1})$, which is also the value of the parameter associated to this edge. Suppose we fix a node $n = a_1 \cdots a_k \in L_k$. Each edge coming out of n gives the probabilities for the variable in the next level L_{k+1} , conditioned on the context $(X_{\pi_1} = a_1, \dots, X_{\pi_k} =$

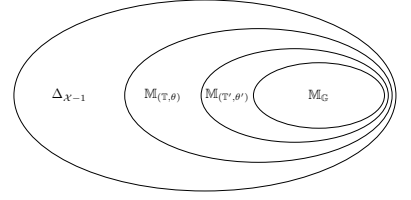


Figure 3.1: Hierarchy of models on \mathcal{X} that we are concerned with, the most general space of models being probability simplex since it contains all distributions on \mathcal{X} , followed by the space of staged tree models $\mathbb{M}_{(\mathbb{T}, \theta)}$, the space of Cstree models $\mathbb{M}_{(\mathbb{T}', \theta')}$ then DAG models \mathbb{M}_G . More general models can explain more datasets whilst simpler models can often be easier to work with.

² The k^{th} level of a rooted tree, L_k , is the set of nodes such that the unique path from each node in L_k to the root consists of k edges.

³ The chain rule in probability states $\mathbb{P}(X_1, \dots, X_p) = \mathbb{P}(X_p | X_{p-1}, \dots, X_1) \mathbb{P}(X_{p-1} | X_{p-2}, \dots, X_1) \cdots \mathbb{P}(X_2 | X_1) \mathbb{P}(X_1)$

a_k). Thus, we can view this node n as containing the distribution $\mathbb{P}(X_{\pi_{k+1}} | X_{\pi_1} = a_1, \dots, X_{\pi_k} = a_k)$. This is an important view which we will make use of when testing for context specific independence in the algorithms throughout this paper. We show an example of a CStree and a staged tree that is not a CStree below.

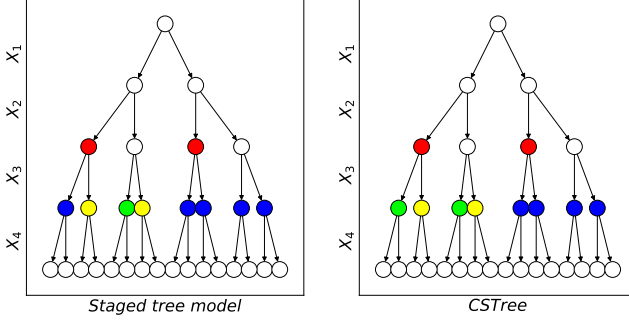


Figure 3.2: Example of a staged tree model that is not a CStree (Left) and a CStree (right) for binary variables X_1, X_2, X_3, X_4 in that causal ordering.

Both staged trees in Figure 3.2 represent 4 binary variables X_1, X_2, X_3, X_4 taking values in $\{0, 1\}$ in that causal order. Suppose each edge to the left corresponds to the outcome 0 and the other corresponds to one. In this case, the left edge coming out of the root represents $\mathbb{P}(X_1 = 0)$ and the right edge coming out of the root represents $\mathbb{P}(X_1 = 1)$. The nodes represent distributions conditioned on the context unique to them. For example, the left most red node in both trees represent $\mathbb{P}(X_3 | X_2 = 0, X_1 = 0)$. The tree on the right is a CStree because each of the nodes in the non-singleton stages, which are represented by a non-white colour, share exactly one fixed context. For example, the stage corresponding to the blue nodes in the tree on the right (the CStree) corresponds to the contexts $(X_1 = 1, X_2 = 0, X_3 = 0)$, $(X_1 = 1, X_2 = 1, X_3 = 1)$, $(X_1 = 1, X_2 = 1, X_3 = 0)$, $(X_1 = 1, X_2 = 0, X_3 = 1)$. The common context for this stage is $(X_1 = 1)$. Meanwhile, for the tree on the left, the stage corresponding to the blue nodes only share the empty context, meaning all nodes in level 3 must correspond to the stage with the empty context for it to be a CStree - this is however not the case since there are nodes in level 3 which correspond to the yellow and green stages, thus not part of the blue stage.

A CStree encodes Context Specific Independence (CSI) relations according to the following lemma ⁴

Lemma 3.5 (CStrees and Context Specific Independence relations)

Let $\mathbb{T} = (\mathbb{V}, \mathbb{E})$ be a CStree with levels $X_1, \dots, X_p \sim L_1, \dots, L_p$ and stages S_1, \dots, S_m . Then for any $\mathbb{P} \in \mathbb{M}_{(\mathbb{T}, \theta)}$ and $S_i \subset L_{k-1}$, \mathbb{P} entails the CSI relation $(X_k \perp\!\!\!\perp_{\mathbb{P}} X_{[k-1] \setminus C_i} | X_{C_i} = x_{C_i})$ where $X_{C_i} = x_{C_i}$ is the context fixed by the stage S_i .

⁴ Eliana Duarte and Liam Solus. Representation of context-specific causal models with observational and interventional data, 2021

The CSI relations from the CStrees look similar to the CI relations from the ordered Markov property. The difference is that the CStree encodes independence of X_k with all the variables preceding it in the causal ordering when conditioned on a context, compared to the ordered Markov property which contains the variables which represent the parents of k in the conditioning set. Thus, CStrees can be thought of as a relaxation of DAG models via a relaxation of the ordered Markov property, where we condition on the more general scenario of contexts, rather than variables.

3.2 Learning CStrees from observed data

Given a system of variables X_1, \dots, X_p , we would first need a causal ordering $\pi_1 \cdots \pi_p$ in order to construct a CStree for these variables. Since CStrees encode CSI relations, they can also encode CI relations, which means we can generate a CStree from a DAG. The following proposition formalizes this notion ⁵.

⁵ Eliana Duarte and Liam Solus. Algebraic geometry of discrete interventional models, 2020

Proposition 3.6 (CStrees corresponding to DAGs) *A compatibly labelled staged tree \mathbb{T} with causal ordering $\pi_1 \cdots \pi_p$, levels L_1, \dots, L_p corresponding to variables $X_{\pi_1}, \dots, X_{\pi_p}$ encodes the same CI relations as some DAG \mathbb{G} if and only if for any topological ordering of \mathbb{G} , $\forall k \in [p-1]$, the level L_k has its nodes partitioned into stages where the context for each stage is an element of the Cartesian product of the parents of $X_{\pi_{k+1}}$ in \mathbb{G} .*

We describe the computational procedure to generate a CStree \mathbb{T} from a DAG \mathbb{G} below, assuming that we are given a causal ordering of \mathbb{G} . ⁶

⁶ PARENTS is a function that takes a graph and a node and returns the parents of that node in the graph; CARTESIANPRODUCT takes a set of variables and returns the cartesian product of these variables i.e. all possible values they can take

Algorithm 1: DAGToCSTREE
Constructing a CStree from a DAG

Input: A DAG G , causal ordering O
Output: CStree T with ordering O and stages S defined by G

```

1  $T \leftarrow$  Empty staged tree with ordering  $O$ ;
2  $S \leftarrow$  Empty List;
3 for  $k$  in  $|O| - 1$  do
4    $v \leftarrow O[k + 1]$  ;
5    $T.add\_level(v)$ ;
6    $pars \leftarrow \text{PARENTS}(G, v)$ ;
7    $contexts \leftarrow \text{CARTESIANPRODUCT}(pars)$ ;
8    $stages\_k \leftarrow$  Empty Dictionary;
9   for  $c$  in  $contexts$  do
10     $stages\_k[c] \leftarrow$  [nodes in level  $k$  where  $c$  is a subcontext];
11  end
12   $S.append(stages\_k)$ ;
13 end
14 return  $T, S$ 

```

Algorithm 1 above does not necessarily need a causal ordering. This is because given a DAG we can perform a topological sort on it to get one, for which efficient algorithms exist ⁷.

Theorem 3.7 *Given variables X_1, \dots, X_p taking values in $\mathcal{X} = \prod_{i=1}^p \mathcal{X}_i$, Algorithm 1 is correct and runs in $\mathcal{O}(d^{2p})$ time and $\mathcal{O}(d^p)$ space where $d = \max_{i \in [p]} |\mathcal{X}_i|$.*

Proof: For correctness, at each level L_k , the non-singleton stages are created for the contexts fixed by the outcomes of the parents of $X_{\pi_{k+1}}$ thus by Proposition 3.6 the tree is still a CStree. Since the staging process at each level only creates non-singleton stages of nodes within that level, and we go over each level except the last level which always contain singleton stages (one for each outcome of \mathcal{X}), the stages S lead to T being a CStree. For time complexity, the worst case scenario is for the fully connected DAG, assuming the ordering $12 \dots p$, node i has $i - 1$ parents. This however results in a CStree with no non-singleton stages. Thus we look at the scenario where node i has $i - 2$ parents. At the level for the variable representing node i , the variable b in Algorithm 1 which is all the elements of the the Cartesian product of values the parents take, has $|\prod_{j=1}^{i-2} \mathcal{X}_j|$ elements. For each element in this Cartesian product which fixes the context for the stage, we have to loop over all nodes in level i and to store the nodes for that stage, and level i has $|\prod_{j=1}^i \mathcal{X}_j|$ nodes. Thus the loop for level i takes $|\prod_{j=3}^{i-2} \mathcal{X}_j| |\prod_{j=1}^i \mathcal{X}_j|$ where the indexing starts at 3 for the first term since the parent sets are non-empty starting from node 3. Since we have p levels, ignoring the first 2 since their variables have no parents, we have

⁷ Robert Endre Tarjan. Edge-disjoint spanning trees and depth-first search. *Acta Informatica*, 6(2):171–185, 1976

$$\sum_{i=3}^p \prod_{j=3}^{i-2} |\mathcal{X}_j| \prod_{k=1}^i |\mathcal{X}_k| < \sum_{i=1}^p \prod_{j=1}^i |\mathcal{X}_j| \prod_{k=1}^i |\mathcal{X}_k| < \sum_{i=1}^p \prod_{j=1}^i d \prod_{k=1}^i d$$

where $d = \max_{i \in [p]} |\mathcal{X}_i|$. This sum then becomes

$$\sum_{i=1}^p d^{2i} = \frac{d^2(d^{2p} - 1)}{d^2 - 1} = \mathcal{O}(d^{2p})$$

For space complexity, in the worst case DAG mentioned, level i which has $\prod_{j=1}^i |\mathcal{X}_j| < d^i$ nodes and the same amount of edges coming in. For storing the stages, the extra information we need to store is the fixed contexts for each stage, and there are $\prod_{j=3}^i |\mathcal{X}_j| < d^i$ stages in level i . Thus the nodes, edges and stages for level i are at most $3d^i$, summing for each level gives

$$\sum_{i=1}^p 3d^i = \frac{3d(3d^p - 1)}{3d - 1} = \mathcal{O}(d^p)$$

We mention the space complexity here to emphasize that it grows exponentially, which is one limitation of this approach. For p binary variables this means a CStree takes $\mathcal{O}(2^p)$ space. This is in comparison to DAGs which in the worst case assuming full connectivity require $\mathcal{O}(p^2)$ space, independent of the state space of the variables.

In order to learn CSI relations, one can now take a CStree from a DAG and perform a statistical test to determine context specific independence relations. Recall that each node in level k represents a probability density of the variable in level $k + 1$ under the context fixed by that node. Thus for each level, we can compare all possible pairs of nodes by taking the samples fixed by the contexts of the pair, and testing whether they are from the same distribution. If so, we assign the same colour to both of them. Then by the CStree property from Definition 3.4 we must have that all nodes in level k which share the same context as that of these 2 nodes must also have the same colour. For example with binary variables if we have 2 nodes representing the outcomes $X_{\{1,2,3,4\}} = 0110, X_{\{1,2,3,4\}} = 0011$ ⁸ and we know they are in the same stage S_i , then the common context for that stage is $X_{\{1,4\}} = 01$, and by the CStree property all nodes in that level with this subcontext belong to the same stage.

⁸ $X_{\{1,2,3,4\}} = 0110$ is shorthand for $(X_1 = 0, X_2 = 1, X_3 = 1, X_4 = 0)$

We now describe the algorithm for learning a CStree.⁹

Algorithm 2: LEARN CSTREE

Learning a CStree with knowledge of causal ordering

Input: CStree T , (possibly empty) stages S , causal ordering O ,
Data matrix D

Output: The CStree T with ordering O and stages S

```

1  $l = 1$ ;
2  $p = |O|$ ;
3 while  $l < p$  do
4    $ns \leftarrow$  [nodes in level  $l$  of  $T$ ];
5    $ps \leftarrow$  [all pairs of nodes in level  $l$ ];
6   for  $(n_1, n_2)$  in  $ps$  do
7     if COLOUR( $n_1$ )=COLOUR( $n_2$ ) then
8       skip
9     else
10       $c \leftarrow$  COMMONCONTEXT( $n_1, n_2$ );
11       $same\_distr =$  TEST( $c, n_1, n_2, D, O[l+1]$ );
12      if  $same\_distr$  then
13         $new\_nodes \leftarrow$  NODESWITHCONTEXT( $ns, c$ );
14         $S \leftarrow$  UPDATESTAGES( $S, c, new\_nodes$ );
15      end
16    end
17  end
18   $l = l + 1$ ;
19 end
20 return  $T, S$ 

```

Algorithm 2 can be sped up by already providing a non-empty CStree containing stages which we may have inferred from a DAG. If one knows the DAG and the true causal ordering of the system they can learn a CStree by using Algorithm 1 followed by Algorithm 2.

Theorem 3.8 *Given variables X_1, \dots, X_p taking values in $\mathcal{X} = \prod_{i=1}^p \mathcal{X}_i$, Algorithm 2 is correct and runs in $\mathcal{O}(d^{2p})$ time assuming constant time for statistical independence testing, where $d = \max_{i \in [p]} |\mathcal{X}_i|$.*

Proof: For correctness, at level loop iteration we compare all pairs of nodes and only update the stages if they do not belong to the same non-singleton stage. In this case if they do belong to the same stage according to the statistical testing, we add exactly the nodes that belong to the stage according to the CStree property. Thus the CStree property is intact throughout the algorithm. For time complexity, using notation from Theorem 3.7, level i has d^i nodes and in the worst case we run statistical testing on all pairs of nodes, of which there are $\binom{d^i}{2} = \frac{d^i!}{2!(d^i-2)!} < \frac{d^{2i}}{2}$, summing for each level

⁹ COLOUR is a function that takes a node and returns the colour of it if it belongs to a non-singleton stage - note here we represent the stage using a colour; COMMONCONTEXT is a function that takes 2 nodes and returns their common context - if one or both of them already belong to a stage, we take this to be the common context between these contexts; TEST is a function that determines whether the distributions corresponding to both of the nodes belong to the same stage or not - this typically involves a statistical test; NODESWITHCONTEXT takes a set of nodes and a context c and returns the nodes which have the c as a subcontext; UPDATESTAGES is a function that updates the stages of the tree with the new nodes.

gives $\sum_{i=1}^p \frac{d^{2i}}{2} = \mathcal{O}(d^{2p})$.

In the general case it is possible that the true causal ordering is unknown. In fact, we need to consider the set of all causal orderings for each DAG in the MEC of the true DAG. Thus we first learn the CPDAG of the true underlying DAG using the PC algorithm and then apply Algorithms 1 and 2.

Algorithm 3: CStreePCALGORITHM

Learning a CStree from observational data

Input: Data matrix D , (optional causal ordering O)

Output: List of CStrees T with minimum number of stages

```

1 CPDAG  $\leftarrow$  PCALGORITHM( $D$ );
2 if  $O$  given then
3    $G \leftarrow g \in \text{CPDAG}$  with ordering  $O$ ;
4    $dags \leftarrow [G]$ ;
5    $orderings \leftarrow [O]$ ;
6 else
7    $dags \leftarrow [g \text{ in CPDAG}]$ ;
8  $min\_stage\_trees \leftarrow []$ ;
9  $min\_stage \leftarrow \infty$ ;
10 for  $G$  in  $dags$  do
11   if  $O$  not given then
12      $orderings \leftarrow \text{ALLTOPOLOGICALSORT}(G)$ ;
13   end
14   for  $O$  in  $orderings$  do
15      $T, S \leftarrow \text{DAGToCSTREE}(G, O)$ ;
16      $T, S \leftarrow \text{LEARNCSTREE}(T, S, O, D)$ ;
17     if  $|S| < min\_stages$  then
18        $min\_stages \leftarrow |S|$ ;
19        $min\_stage\_trees \leftarrow [(T, S)]$ ;
20     end
21     if  $|S| = min\_stages$  then
22        $min\_stage\_trees.append((T, S))$ ;
23     end
24   end
25 end
26 return  $min\_stage\_trees$ 

```

We consider all topological orderings of all Markov equivalent DAGs learnt by the PC algorithm because we might not be able to encode some context specific information otherwise. We show an example of this in the next section after introducing minimal context DAGs.

Algorithm 3 considers many possible candidate CStree models, thus we have to pick the best model with respect to some criterion. There are however many instances where we could know the causal ordering apriori ¹⁰, for example a temporal relation between nodes known through physical laws. In this case we can either start statistical testing from an empty tree, or apply the PC algorithm to the data and find a DAG in the Markov Equivalence class with the known ordering and then run the extra testing.

Unlike DAGs, as the number of variables increase it gets progressively harder to visually understand the learnt causal structure by just looking at the learnt CStree.

3.3 Understanding high-dimensional CStrees

From a pragmatic perspective the aim of this section is to introduce the notion of Minimal Context (MC) DAGs which can help visualize CStrees with more variables and the context specific information they encode. On a theoretical note, this work has led to the generalization of Theorem 2.9 to define a characterization of Markov Equivalence for CStrees ¹¹. We start by first describing the procedure ¹² to generate the CSI relations from a CStree and its stages, which uses Lemma 3.5

Algorithm 4: GENERATECSIRELATIONS

Generate the CSI Relations from the CStree

Input: CStree T , its stages S and its causal ordering O

Output: Set of CSI Relations J encoded in the CStree

```

1  $l = 1$ ;
2  $p = |O|$ ;
3  $J = []$ ;
4 while  $l < p$  do
5    $S_l \leftarrow \text{STAGESINLEVEL}(S, l)$ ;
6   for  $s$  in  $S_l$  do
7      $c \leftarrow \text{CONTEXTOfSTAGE}(s)$ ;
8      $v_c \leftarrow \text{VARIABLESOfCONTEXT}(c)$ ;
9      $v_o \leftarrow O[1 : l - 1] \setminus v_c$ ;
10     $J.append((X_{O[l+1]} \perp\!\!\!\perp X_{v_o} \mid c))$ 
11  end
12 end

```

Theorem 3.9 Given a CStree \mathbb{T} , Algorithm 4 is correct and returns the CSI relations encoded in \mathbb{T} in $\mathcal{O}(pd^{2p})$ time.

Proof: Correctness follows directly from Lemma 3.5 since at each loop we

¹⁰ Peter Thwaites, Jim Q. Smith, and Eva Riccomagno. Causal analysis with chain event graphs. *Artificial Intelligence*, 174(12):889–909, 2010; and Tomi Silander and Tze-Yun Leong. A dynamic programming algorithm for learning chain event graphs. In *International Conference on Discovery Science*, pages 201–216. Springer, 2013

¹¹ Eliana Duarte and Liam Solus. Representation of context-specific causal models with observational and interventional data, 2021

¹² STAGESINLEVEL takes a set of stages and a level and returns the stages in that level; CONTEXTOfSTAGE takes a stage and returns the common context of that stage; VARIABLESOfCONTEXT takes a context and returns the variables in it

add exactly the CSI relations mentioned in the lemma, and we do this for all levels thus include all stages of the CStree. For time complexity, for each level we first get the stages associated with it, which can be done in constant time if we store this information. We know from the Proof of Theorem 3.7 that the number of stages in level i is bounded above by d^i , and for each stage in that level we get the context of the stage, which can be done as a constant lookup operation, and we get relevant variables in Lines 8,9 which is bounded above by $2p$. Thus adding this for all the levels give $\sum_{i=1}^p 2pd^i = \mathcal{O}(pd^{2p})$.

In practice a slightly modified version of Algorithm 4 can be placed a subroutine in the previous algorithms right before moving onto the next level.

From Lemma 3.5 we know any distribution in the CStree model $\mathbb{M}_{(\mathbb{T},\theta)}$ encodes the CSI relations of the given form, there could be more CSI relations satisfied by every distribution in $\mathbb{M}_{(\mathbb{T},\theta)}$, which is similar to how a DAG model encodes the CI relations \mathbb{J}_1 implied by the local Markov property, which are captured by the CI relations \mathbb{J}_2 from the global Markov property (i.e. from the d-separations), and $\mathbb{J}_1 \subset \mathbb{J}_2$.

The complete set of all CSI relations satisfied by each distribution $\mathbb{P} \in \mathbb{M}_{(\mathbb{T},\theta)}$ includes the CSI relations recovered from Algorithm 4, and also include those implied by the successive application of the Context Specific Conditional Independence axioms to generate further CSI relations. The axioms are as follows

1. Symmetry, If $(X_A \perp\!\!\!\perp X_B \mid X_C = x_c) \in \mathbb{J}$ then $(X_B \perp\!\!\!\perp X_A \mid X_C = x_c) \in \mathbb{J}$
2. Decomposition, If $(X_A \perp\!\!\!\perp X_{B \cup D} \mid X_S, X_C = x_c) \in \mathbb{J}$ then $(X_A \perp\!\!\!\perp X_B \mid X_S, X_C = x_c) \in \mathbb{J}$
3. Weak union, If $(X_A \perp\!\!\!\perp X_{B \cup D} \mid X_S, X_C = x_c) \in \mathbb{J}$ then $(X_A \perp\!\!\!\perp X_B \mid X_{S \cup D}, X_C = x_c) \in \mathbb{J}$
4. Contraction, If $(X_A \perp\!\!\!\perp X_B \mid X_{S \cup D}, X_C = x_c) \in \mathbb{J}$ and $(X_A \perp\!\!\!\perp X_D \mid X_S, X_C = x_c) \in \mathbb{J}$ then $(X_A \perp\!\!\!\perp X_{B \cup D} \mid X_S, X_C = x_c) \in \mathbb{J}$
5. Intersection, If $(X_A \perp\!\!\!\perp X_B \mid X_{S \cup D}, X_C = x_c) \in \mathbb{J}$ and $(X_A \perp\!\!\!\perp X_B \mid X_{B \cup D}, X_C = x_c) \in \mathbb{J}$ then $(X_A \perp\!\!\!\perp X_{B \cup S} \mid X_D, X_C = x_c) \in \mathbb{J}$
6. Specialization, If $(X_A \perp\!\!\!\perp X_B \mid X_S, X_C = x_c) \in \mathbb{J}$ and $T \subset S, x_T \in \mathcal{X}_T$ then $(X_A \perp\!\!\!\perp X_B \mid X_{S \setminus T}, X_{T \cup C} = x_{T \cup C}) \in \mathbb{J}$
7. Absorption, If $(X_A \perp\!\!\!\perp X_B \mid X_S, X_C = x_c) \in \mathbb{J}$ and $\exists T \subset C$ such that $\forall x_T \in \mathcal{X}_T$ we have $(X_A \perp\!\!\!\perp X_B \mid X_S, X_{C \setminus T} = x_{C \setminus T}, X_T = x_T) \in \mathbb{J}$ then $(X_A \perp\!\!\!\perp X_B \mid X_{S \cup T}, X_{C \setminus T} = x_{C \setminus T}) \in \mathbb{J}$

Given a Context Specific Conditional Independence model \mathbb{J} , the successive application of the axioms above results in the Context Specific closure of the model, denoted $\bar{\mathbb{J}}$.

The Absorption axiom helps us to get a representation of the CStree as a sequence of DAGs. For this we need the definition of minimal contexts ¹³.

Definition 3.10 (Minimal contexts) *Given a set Context Specific Independence model \mathbb{J} , we say that $X_C = x_C$ is a minimal context if we have $(X_A \perp\!\!\!\perp X_B \mid X_S, X_C = x_C) \in \mathbb{J}$ and there is no non-empty subset $T \subset C$ such that $(X_A \perp\!\!\!\perp X_B \mid X_{S \cup T}, X_{C \setminus T} = x_{C \setminus T}) \in \mathbb{J}$.*

Intuitively, the minimal contexts are the smallest contexts that get left behind after repeated application of the Absorption axiom. By the Specialization axiom, given a minimal context $(X_C = x_C)$ we can recover all the CI relations implied by the CStree.

Minimal contexts and the application of the graphoid axioms are key to represent the CStree as a sequence of DAGs. For now, we denote $\mathbf{C}(\mathbb{T})$ as the set of all minimal contexts of a CStree \mathbb{T} , for each minimal context $(X_C = x_C)_i \in \mathbf{C}(\mathbb{T})$, there exists a DAG that encodes the CI relations that hold under this context. We denote $\mathbf{G}(\mathbb{T}) := \{\mathbf{G}_{X_C=x_C}\}_{X_C=x_C \in \mathbf{C}(\mathbb{T})}$ as the set of all minimal context DAGs of \mathbb{T} , and define the procedure to generate them later on.

We give some examples of minimal contexts below.

Example 3.11 *Let X_1, X_2, X_3, X_4, X_5 be binary variables taking values in $\{0, 1\}$. If we have just the CSI relations $(X_5 \perp\!\!\!\perp X_4 \mid X_{\{1,2,3\}} = 000)$ and $(X_5 \perp\!\!\!\perp X_4 \mid X_{\{1,2,3\}} = 100)$ they get absorbed into $(X_5 \perp\!\!\!\perp X_4 \mid X_{\{1\}}, X_{\{2,3\}} = 00)$ leaving the minimal context $X_{\{2,3\}} = 00$, or simply $(X_2 = 0, X_3 = 0)$.*

Example 3.12 *Let X_1, X_2, X_3, X_4 be binary variables taking values in $\{0, 1\}$. If we have the CSI relations $(X_4 \perp\!\!\!\perp X_2 \mid X_{\{1,3\}} = 00)$, $(X_4 \perp\!\!\!\perp X_2 \mid X_{\{1,3\}} = 01)$, $(X_4 \perp\!\!\!\perp X_2 \mid X_{\{1,3\}} = 10)$, $(X_3 \perp\!\!\!\perp X_1 \mid X_2 = 1)$ then they absorb to give the equivalent CSI relations $(X_4 \perp\!\!\!\perp X_2 \mid X_1, X_3 = 0)$, $(X_4 \perp\!\!\!\perp X_2 \mid X_3, X_1 = 0)$, $(X_3 \perp\!\!\!\perp X_1 \mid X_2 = 0)$ thus the minimal contexts are $(X_1 = 0)$, $(X_2 = 0)$, $(X_3 = 0)$*

Example 3.11 shows that we can get absorption from CSI relations we get from different levels, while 3.12 shows that we have to check at least all possible pairs of CSI relations to get the minimal contexts. It is also possible to get the empty context as the minimal context, which happens when the all CSI relations $(X_A \perp\!\!\!\perp X_B \mid X_S, X_C = x_C)$ hold for all outcomes $x_C \in \mathcal{X}_C$ in which case the absorption axiom gives the CI relation $(X_A \perp\!\!\!\perp X_B \mid X_S, X_C)$.

¹³ Eliana Duarte and Liam Solus. Representation of context-specific causal models with observational and interventional data, 2021

From a computational perspective, given all CSI relations involving sets of variables X_A, X_B, X_S i.e. those of the form $(X_A \perp\!\!\!\perp X_B \mid X_S, X_C = x_C)$, we want to find the largest set $T \subset C$ such $(X_A \perp\!\!\!\perp X_B \mid X_S, X_T = x_T, x_{C \setminus T} = x_{C \setminus T})$ is also in the CSI relations.¹⁴ This can be thought of as decomposing C into sets $C \setminus T$ and T such that Definition 3.10 holds. To find the largest such T , we can perform a binary search on the size of T , which can range from 0 to $|C|$. Algorithm 5 describes our approach, which makes use of Algorithm 6 as a subroutine to find the largest possible T as in the definition of the absorption axiom.¹⁵

Theorem 3.13 *Given a set of Context Specific Independence relations \mathbb{J} , Algorithm 5 returns the set of minimal contexts.*

Proof: We work in sets of CSI relations where the variables A, B and the conditioning set S remain the same, which we denote $\mathbb{J}_{A,B,S}$. For each CSI relation $r \in \mathbb{J}_{A,B,S}$ which has a context with context variables C , Algorithm 6 finds the largest subset T of C such that $\mathbb{J}_{A,B,S}$ contains relations with contexts containing all outcomes of T as a subcontext. The possible sizes of T have the property that if we have a set of size k_1 that contains all of its outcomes in $\mathbb{J}_{A,B,S}$ then all sets of size $0, \dots, k_1 - 1$ also have all their outcomes in $\mathbb{J}_{A,B,S}$ as subcontexts meaning in order to maximize this size we then have to search over the sizes $k_1, \dots, |C|$. Similarly, if we find a set of size k_2 which does not contain a set with all possible sizes for T inside $\mathbb{J}_{A,B,S}$, we know this size will be in $0, \dots, k_2 - 1$. Noting that the set with size 0 i.e. the empty set is always a subcontext, and since the algorithm searches over all possible sizes for T (which can only range from 0 to $|C|$ since it is a subset of C) we get the largest set T such that the absorption axiom is satisfied. Once we find the size of T , we check for T itself. For each possible T_i with the size we just found, this is done by decomposing the context variable of r , into T_i and $C \setminus T_i$ and ensuring all outcomes of T_i have a corresponding CI relation in $\mathbb{J}_{A,B,S}$ and the corresponding possible minimal context after this decomposition contains the context $X_{C \setminus T_i} = x_{C \setminus T_i}$ as a subcontext, which is what the absorption axiom asks for.

A crude time complexity approximation is as follows. We start by computing how long it would take to find the optimal size for T given a CSI relation r and a set of CSI relations \mathbb{J} which share the same variables and conditioning set. In the worst case, the possible variables for T can be all the variables, which we denote p . Let d denote the maximum number of outcomes over all variables. The first midpoint is then $p/2$, which gives a total of $\binom{p}{2}$ possible subsets of p with this size. For each of set, we generate the Cartesian product, which has at most $d^{\frac{p}{2}}$ elements. Generalizing this, at iteration i , we will look over $\binom{p}{2^i}$ possible subsets and corresponding $d^{\frac{p}{2^i}}$

¹⁴ If we are to use the CSI relations extracted from Algorithm 4 before getting their Context Specific closure they would be of the form $(X_i \perp\!\!\!\perp X_B \mid X_C = x_C)$ where $i \in [p]$ and $B, C \subset [p - 1]$

¹⁵ Each break statement leaves the first loop it meets. `RELSWITH-TRIPLE(A, B, S, J)` returns the CSI relations with the variables A, B and conditioning set S in J . `VARIABLESOFCONTEXT` takes a context and returns the variables involved in it. `CONTAINED` takes a possible outcome x_T , a list of contexts, and a context $X_L = x_L$ (which can be a possible minimal context) and returns true if there is a context $X_C = x_C$ in the list of contexts such that it can be decomposed into $X_{C \setminus T} = x_{C \setminus T}, X_T = x_T$ and $X_L = x_L$ is a subcontext of $X_C = x_C$. `UPDATEMINIMALCONTEXTS` is a function that updates the dictionary which stores the minimal contexts alongside all the CSI relations in J that include it as a subcontext/

Algorithm 5: GENERATEMINIMALCONTEXTS

Generate the Minimal Contexts from a set of CSI relations

Input: Set of CSI Relations J **Output:** Set of minimal contexts MCs from J

```

1  $set\_triples \leftarrow \text{TRIPLESOFSETS}(J);$ 
2  $MCs \leftarrow \text{EMPTYDICTIONARY};$ 
3 for  $A, B, S$  in  $set\_triples$  do
4    $rels \leftarrow \text{RELSWITHTRIPLE}(A, B, S, J);$ 
5    $cs \leftarrow \text{CONTEXTSOFRELS}(rels);$ 
6   for  $rel$  in  $rels$  do
7      $rel\_context \leftarrow \text{CONTEXTOFRELS}([rel]);$ 
8      $C \leftarrow \text{VARIABLESOFCONTEXT}(rel\_context);$ 
9      $T\_sizes \leftarrow [0, \dots, |C|];$ 
10     $T\_size \leftarrow \text{RECURSIVESTRUCT}(T\_sizes, C, cs);$ 
11    if  $T\_size=1$  then
12       $T = \{\};$ 
13    else if  $T\_size = |C|$  then
14       $T = C;$ 
15    else
16       $possible\_Ts \leftarrow [\text{subsets of } C \text{ of size } T\_size];$ 
17      for  $possible\_T$  in  $possible\_Ts$  do
18         $all\_outcomes \leftarrow \text{CARTESIANPRODUCT}(possible\_T);$ 
19         $possible\_mc \leftarrow [(var, val) \text{ for } (var, val) \text{ in}$ 
20           $rel\_context \text{ if } var \text{ not in } possible\_T];$ 
21         $outcome\_count = 0;$ 
22        for  $outcome$  in  $all\_outcomes$  do
23           $contained \leftarrow$ 
24             $\text{CONTAINED}(outcome, cs, possible\_mc);$ 
25          if  $contained = \text{True}$  then
26             $outcome\_count = outcome\_count + 1;$ 
27          end
28          if  $outcome\_count = |all\_outcomes|$  then
29             $T = possible\_T$ 
30          end
31        end
32      end
33    end
34  end
  
```

Algorithm 6: RECURSIVESHARCHT

Find the largest possible size of the set T for which the absorption axiom can be applied.

Input: Possible sizes for T denoted T_sizes , possible variables for T denoted T' , set of contexts shared by the same triple A, B, S , denoted cs

Output: Size of the largest set T such that all outcomes under T are contained in cs

```

1 if  $|T\_sizes| = 1$  then
2   return  $T\_sizes[0]$ 
3 else
4    $mid \leftarrow |T\_sizes|/2$ ;
5    $candidate\_size \leftarrow T\_sizes[mid]$ ;
6    $possible\_Ts \leftarrow$  [subsets of  $C$  of size  $candidate\_size$ ];
7    $not\_contained = 0$ ;
8   for  $T$  in  $possible\_Ts$  do
9      $all\_outcomes \leftarrow$  CARTESIANPRODUCT( $T$ )
10     $outcome\_count = 0$ ;
11    for  $outcome$  in  $all\_outcomes$  do
12       $contained \leftarrow$  CONTAINED( $outcome, cs, \{\}$ );
13      if  $contained = \text{True}$  then
14         $outcome\_count = outcome\_count + 1$ ;
15      end
16      if  $contained = \text{False}$  then
17         $not\_contained = not\_contained + 1$ ;
18        break
19      end
20      if  $outcome\_count = |all\_outcomes|$  then
21        RECURSIVESHARCHT( $T\_sizes[mid:end], C, cs$ )
22      end
23    end
24    if  $not\_contained = |possible\_Ts|$  then
25      RECURSIVESHARCHT( $T\_sizes[start:mid], C, cs$ )
26    end
27  end

```

outcomes. We do this at most $\log p$ times, giving the time complexity to be

$$\sum_{i=1}^{\log p} \binom{p}{\frac{p}{2^i}} d^{\frac{p}{2^i}}$$

Once we find the optimal size for T , we then find the set T itself, which involves generating the Cartesian product over all possible subsets of the size itself. This is done for all CSI relations of the form $(X_A \perp\!\!\!\perp X_B | X_S, X_{C_i} = x_{C_i})$. Since we look at all pairs of sets A, B , and given that they are in general, subsets of $[p]$, we have $\binom{2^p}{2}$ possibilities for the pairs of sets A, B . In practice, the number of CSI relations we have to look at can be reduced, for example the Symmetry axiom allows us to half the number of pairs we need to check. This motivates us to the pairwise case, where we focus on Context Specific Independence relations of the form between 2 variables X_i, X_j rather than two sets of variables X_A, X_B . This drops the number of possible pairs to $\binom{p}{2}$. We note that however that even though we fix pairs of variables X_i, X_j now, we could still have the conditioning set S to be any subset of size $p - 2$ which does not include the nodes corresponding to X_i, X_j . This leads us to the definition of pairwise minimal contexts.

Definition 3.14 (Pairwise minimal contexts) *Given a set Context Specific Independence model \mathbb{J} , we say that $X_C = x_C$ is a pairwise minimal context if we have $(X_i \perp\!\!\!\perp X_j | X_S, X_C = x_C) \in \mathbb{J}$ and there is no non-empty subset $T \subset C$ such that $(X_i \perp\!\!\!\perp X_j | X_{S \cup T}, X_C = x_C) \in \mathbb{J}$.*

We can make use of the Context Specific Conditional independence axioms to generate pairwise CSI relations to get the pairwise minimal contexts.

The introduction of pairwise relations motivates the inclusion of the following axiom.

8. Composition, If $(X_A \perp\!\!\!\perp X_B | X_S, X_C = x_C)$ and $(X_A \perp\!\!\!\perp X_D | X_S, X_C = x_C)$ then $(X_A \perp\!\!\!\perp X_{B \cup D} | X_S, X_C = x_C)$

The composition axiom allows us to go from pairwise relations between variables to the more general pairwise relations between sets of variables. We call a Context Specific Conditional Independence model satisfying the additional composition axiom a Context Specific compositional graphoid, which generalizes the notion of compositional graphoids¹⁶. An important question is whether the Context Specific closure of CSI relations from a CStree form a compositional context specific graphoid.

¹⁶ Kayvan Sadeghi and Steffen Lauritzen. Markov Properties for Mixed graphs. *Bernoulli*, 20(2):676 – 696, 2014

Pairwise minimal contexts are always minimal contexts however but we could have minimal contexts that are not pairwise minimal contexts. If all CStrees have an associated compositional context specific independence model, the pairwise relations may indeed be all that we need to get the set of complete set of minimal contexts. We leave this as an open question and use pairwise minimal contexts to offer a visualization of CSI relations in the CStree which may potentially be incomplete.

We now have the machinery to visualize higher dimensional CStrees. We start by generating the CSI relations from the trees. Then get the Context Specific closure, followed by the minimal contexts. Once we have the minimal contexts, the CI relations for each minimal context define a graphoid. This is called a Minimal Context DAG. Thus, the CStree can be represented as a sequence of DAGs for each minimal context. This representation is a consequence of the following theorem ¹⁷.

Theorem 3.15 (Markov theorem for CStrees) *Given a CStree \mathbb{T} , with levels $L_1, \dots, L_p \sim X_1, \dots, X_p$, minimal contexts $\mathcal{C}(\mathbb{T})$ and $\mathbb{P} \in \Delta_{|\mathcal{X}|-1}$. The following are equivalent.*

- \mathbb{P} factorizes according to \mathbb{T} .
- \mathbb{P} is Markov to $\mathcal{G}(\mathbb{T})$.
- $\forall X_C = x_C \in \mathcal{C}(\mathbb{T})$ we have

$$\mathbb{P}(X_{[p] \setminus C} \mid X_C = x_C) = \prod_{k \in [p] \setminus C} \mathbb{P}(X_k \mid X_{PA_{\mathcal{G}_{X_C=x_C}}(k)}, X_C = x_C).$$

This DAG is also called the minimal I-MAP ¹⁸, which we can recover with the following procedure ¹⁹.

Theorem 3.16 *Given variables X_1, \dots, X_p , a set of minimal contexts \mathcal{C} and for each minimal context C the conditional independence relations \mathbb{J}_C that hold under this minimal context, Algorithm 7 is correct and runs in $\mathcal{O}(p^2|\mathcal{C}||\mathbb{J}|)$ time.*

We end this section by stating the definition of faithfulness for CStrees which allows us to explain why we consider all Markov equivalent DAGs when learning a CStree, followed by the main theorem for Algorithm 3

Definition 3.17 (Faithfulness for CStrees) *A distribution \mathbb{P} is faithful to a CStree \mathbb{T} if it entails exactly the CSI relations encoded by the set of minimal context DAGs $\mathcal{G}(\mathbb{T})$.*

To see why we need to consider all Markov equivalent DAGs in Algorithm 3, suppose we have the following case with 4 binary

¹⁷ Eliana Duarte and Liam Solus. Representation of context-specific causal models with observational and interventional data, 2021

¹⁸ Thomas Verma and Judea Pearl. *Causal Networks: Semantics and Expressiveness. "Graphoids: A Computer Representation for Dependencies and Relevance in Automated Reasoning (Computer Information Science)".*, pages 69–76. Uncertainty in Artificial Intelligence. Elsevier, 1990

¹⁹ L Solus, Y Wang, and C Uhler. Consistency Guarantees for Greedy Permutation-Based Causal Inference Algorithms. *Biometrika*, 01 2021. asaa104

Algorithm 7: GENERATEMINCONTEXTDAGS

 Generating minimal context DAGs

Input: Causal ordering O , Minimal Contexts MCs as a dictionary with minimal contexts as keys and the CI relations under the minimal context as values

Output: List of minimal contexts with their minimal context DAGs $MCDAGS$

```

1   $MCDAGS \leftarrow []$ ;
2  for  $MC, Ci\_Rels$  in  $MCs$  do
3     $G \leftarrow$  Empty Graph;
4     $nodes \leftarrow O \setminus \text{VARIABLESOFCONTEXT}(MC)$ ;
5    for  $i$  in  $[1, \dots, |nodes|]$  do
6      for  $j$  in  $[i + 1, \dots, |nodes|]$  do
7         $\pi_i \leftarrow nodes[i]$ ;
8         $\pi_j \leftarrow nodes[j]$ ;
9         $G.add\_edge(\pi_i, \pi_j)$ ;
10        $S = O[1 : j - 1] \setminus \text{VARIABLESOFCONTEXT}(MC) \setminus \{\pi_i\}$  ;
11       if  $(X_{\pi_i} \perp\!\!\!\perp X_{\pi_j} \mid X_S) \in Ci\_Rels$  then
12          $G.remove\_edge(\pi_i, \pi_j)$ 
13       end
14     end
15   end
16    $MCDAGS.add((MC, G))$ ;
17 end
18 return  $MCDAGS$ 

```

variables.

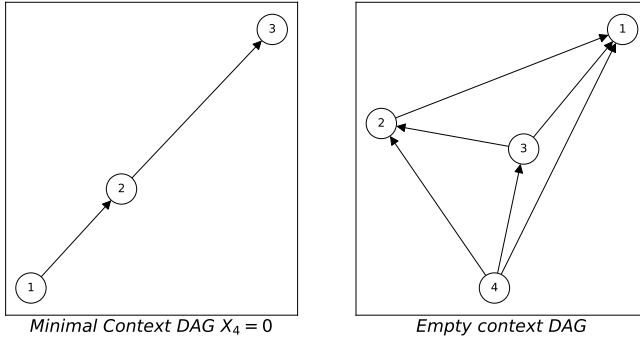


Figure 3.3: Example on why to consider all topological ordering of all Markov equivalent DAGs when learning a CStree

Let \mathbb{P} be the data generating distribution and suppose it is faithful to the CStree with the context graph shown above for the context $X_4 = 0$. The empty context DAG could possibly be the fully connected shown above. If we do happen to learn this DAG from the PC algorithm step, it only has one topological ordering which is 4321, and there is no way to join the stages for the empty context graph which this order to encode $(X_1 \perp\!\!\!\perp X_3 \mid X_4 = 0)$, in which case we do not learn the true model.

Theorem 3.18 *Given variables X_1, \dots, X_p , assuming that the data generating distribution \mathbb{P} is faithful to some unknown CStree \mathbb{T} with a known causal ordering $\pi = \pi_1 \cdots \pi_p$, Algorithm 3 is consistent, i.e. it recovers \mathbb{T} as the number of samples $n \rightarrow \infty$.*

Proof: Since \mathbb{P} is faithful to \mathbb{T} it entails the CSI relations encoded by minimal context DAGs $\mathbb{G}_{\mathbb{T}}$. By Theorem 3.15, this is equivalent \mathbb{P} factorizing according to \mathbb{T} . Faithfulness also implies that if two nodes are in different stages, they must have different labels, since otherwise we would have a CSI relation not entailed by $\mathbb{G}_{\mathbb{T}}$. Thus in the limit of large data, we can differentiate between different stages in \mathbb{T} , allowing us to recover it up to Markov equivalence.

3.4 Model selection for CStrees

In Algorithm 3 we select the CStree (or CStrees if there is more than one) with the causal ordering corresponding to the fewest stages, which is motivated by the principle of Occams razor²⁰ since few stages imply lower model complexity and all else being, the simplest model is the best model. The original authors also present the **Bayesian Information Criterion (BIC)** for CStrees alongside a proof that it is locally consistent for CStrees, meaning it can be

²⁰ Judea Pearl. *Causality*. Cambridge University Press, 2009

applicable to greedy search methods for CStrees. The BIC depends on two terms, the likelihood of the data, which assesses the quality of the model to explain the observed data, and a complexity term that depends on the amount of observed data and the free parameters of the model, with the idea being that models with more parameters should be penalized. Importantly, this helps against overfitting since one can simply add more parameters to the model to maximize the likelihood.

For a general model with random variables X_1, \dots, X_p the likelihood of observing a sample (x_1, \dots, x_p) can be described as below ²¹.

$$\begin{aligned} \mathbb{P}(X_{\{1, \dots, p\}} = x_{\{1, \dots, p\}}) &= \mathbb{P}(X_i = x_i | X_{\{i-1, \dots, 1\}} = x_{\{i-1, \dots, 1\}}) \\ &\quad \cdots \mathbb{P}(X_2 = x_2 | X_1 = x_1) \mathbb{P}(X_1 = x_1) \end{aligned}$$

In the CStree model, independence is implied from the staging of the tree, and each distribution in the product above depends on the fixed context of the node in the CStree representing that distribution. Denoting C_i to be the context variables of the context of the node $x_{1 \dots i}$, this gives the following.

$$\mathbb{P}(X_{\{1, \dots, p\}} = x_{\{1, \dots, p\}}) = \prod_{i=1}^p \mathbb{P}(X_i = x_i | X_{C_i} = x_{C_i})$$

In order to get these values from the data, we denote \mathbb{U} to be **contingency table** of the data, which a multi-dimensional array ²² with dimensions $d_1 \times \dots \times d_p$, where d_i is the number of possible outcomes for variable X_i . Given a sample (x_1, \dots, x_p) , the value in $\mathbb{U}[(x_1, \dots, x_p)]$ is the number of times we have this sample in the dataset. The marginalized contingency table \mathbb{U}_C represents the table after summing over the axes which are not in C . An important case is when C is the empty set, which results in the marginal table to be 1 dimensional scalar and is the total number of samples in the dataset - when having a stage with an empty context this reflects the conditional distribution for the corresponding variable to simply be the ratio of its outcomes in the dataset. This allows the following compact representation for the likelihood ²³ for a tree with levels $(L_1, \dots, L_p) \sim (X_1, \dots, X_p)$.

$$\mathbb{P}(X_{\{1, \dots, p\}} = x_{\{1, \dots, p\}}) = \prod_{k=1}^p \frac{\mathbb{U}_{C \cup k}[x_{C \cup k}]}{\mathbb{U}_C[x_C]}$$

²¹ Here we compactify the notation further - $X_{\{1, \dots, p\}} = x_{\{1, \dots, p\}}$ is the same as $X_{\{1, \dots, p\}} = x_1 \cdots x_p$ which is $X_1 = x_1, \dots, X_p = x_p$

²² This is the same as tensors in the context of computer science

²³ Eliana Duarte and Liam Solus. Representation of context-specific causal models with observational and interventional data, 2021

Given n samples arranged into a $n \times p$ array \mathbb{D} , and under the assumption that they are independent, the total likelihood is simply the product over all samples in \mathbb{D} .

The free parameters is $d = \sum_{k=1}^{p-1} (|\mathcal{X}_k| - 1)S_k$ where S_k is the number of distinct stages in level k ²⁴. The BIC score for a CStree \mathbb{T} with observed data \mathbb{D} is then

$$\text{BIC}(\mathbb{T}, \mathbb{D}) = \log \mathbb{P}(\mathbb{D} \mid \mathbb{T}) - \frac{d}{2} \log(n)$$

²⁴ Eliana Duarte and Liam Solus. Representation of context-specific causal models with observational and interventional data, 2021

4

Experiments

The aim of this section is to run test the theory presented so far in both synthetic and real world data. Namely, for all the datasets we use, we ask ourselves the following.

1. Given a DAG and a causal ordering, is there a difference between learning the context specific independence relations after encoding the CI relations from the DAG into CStree, or without encoding any of these CI relations?
2. Is the CStree property assumption reasonable in practice?
3. How sensitive are the results on the method used to determine whether nodes belong to the same stage?
4. How are the results in comparison to CStrees generated from DAGs learnt using other causal discovery algorithms?
5. How well do our learnt CStrees compared to staged tree models learnt using different algorithms?

Throughout this section we omit the final layer of nodes for all the CStrees since they always belong to the singleton stage.

4.1 Encoding DAGs as CStrees

We first run a sanity check on generating the CStrees from DAGs in the 2 extreme cases, when we have no connectivity and full connectivity.

The results are as expected since there are no conditional independence relations encoded in the fully connected model, and no context specific independence statements since they are unable to encode them, meaning no 2 nodes share a fixed context and all stages have a

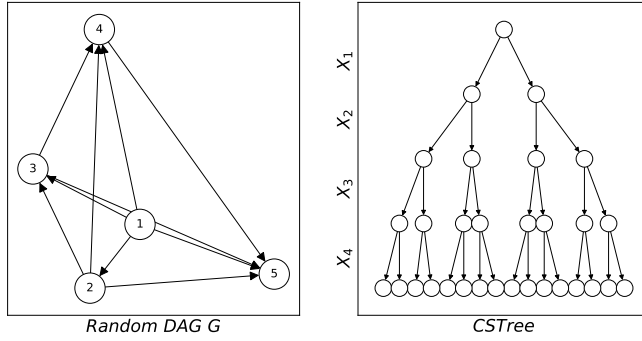


Figure 4.1: Generating the CStrees for a fully connected DAG using the ordering 1234 with Algorithm 1

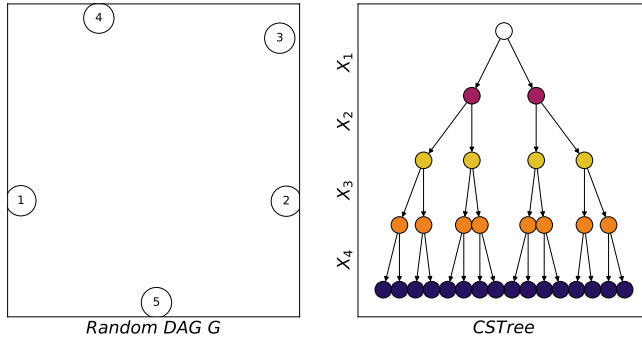


Figure 4.2: Generating the CStrees for DAG with no edges using the ordering 1234 with Algorithm 1

single element. For the DAG with no edges, each variable is independent of all others, and each node has no parents. Thus by Proposition ?? all levels share the empty context as the common context.

4.2 Recovering the empty context DAG

We start by generating random DAGs and generating the CStrees as per Algorithm 1. This is done by first choosing the number of variables (nodes in the DAG) p , and then choosing an edge with probability $p_{edge} \in (0, 1)$, and then keeping edge (u, v) if and only if $u < v$. The causal ordering is chosen as $12 \cdots p$. We show the generated CSTree below, and recover the original DAG as a minimal context DAG with the empty context with Algorithm 7.

The following table summarizes the empirical space time complexity of generating a CSTree from certain random DAGs generated from the aforementioned procedure.

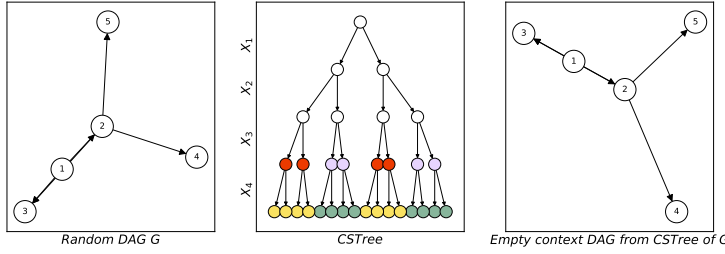


Figure 4.3: Generating CStrees from random DAGs and recovering the original random DAG from the CStree with ordering 1234 using Algorithm 7

p_{edge}	DAG Nodes	Space (GB)	Time (s)
0.2	20	1.203	14
0.2	21	2.422	43
0.2	22	4.875	74
0.2	23	9.812	240

Table 4.1: Statistics from generating CStrees for random DAGs, Space refers to the amount of RAM occupied by the CStree structure, and time refers to the amount of time taken to encode the CI relations in the DAG to the CStree.

4.3 Synthetic data

Next we generate a random DAG G using the procedure above, and then generate samples from this DAG, where each variable is binary valued. We then discard the DAG and apply Algorithm 3 on the generated samples. For nodes which are adjacent to another node, we sample the data by constructing a conditional probability distribution table by taking the parents of the nodes, and creating one row for each possible outcome of the parents. In this binary case, if a variable X_i has d parents $X_{PA_{X_1}}, \dots, X_{PA_{X_d}}$, the corresponding table has 2^d rows. The 2 columns represent the outcome 0 and 1 respectively. The probability values are filled by taking a row corresponding to some outcome of the parents, and assigning the probability of $X_i = 0$ under this outcome to be distributed uniformly between $[0.01, 0.2]$ or $[0.8, 0.99]$ each with a probability 0.5. The probability for $X_i = 1$ under this outcome is then simply computed such that the probabilities sum to 1. For the variables corresponding to nodes with no edges, we give them a Bernoulli distribution with a parameter chosen uniformly between $[0, 1]$ for each such variable.

We start by generating 1000 samples on 4 binary variables using the above procedure. We then learn the CStree and recover the minimal context DAGs. We choose 4 variables due to limitations of applying the context specific graphoid axioms.

The DAG and conditional probability table is as follows.

The above procedure could produce non empty minimal contexts for several reasons, for example due to small sample sizes, or errors in the context specific independence testing procedure.

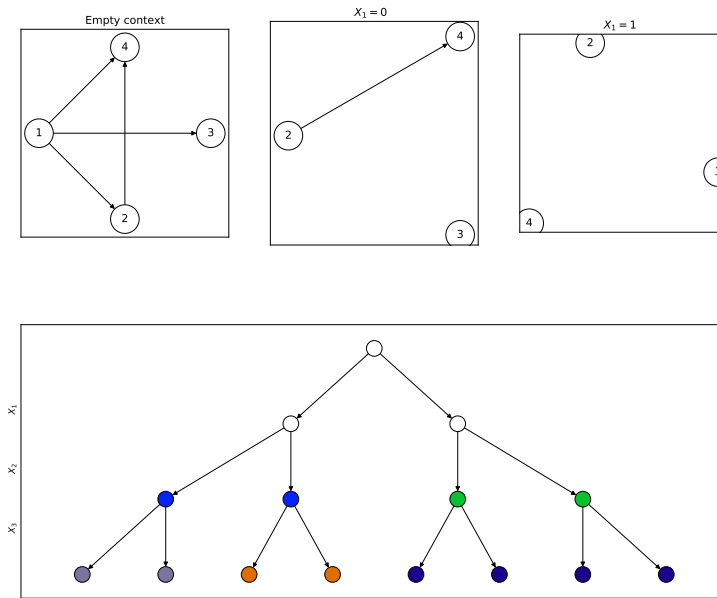


Figure 4.4: Generating CStrees from datasets generated by random DAGs and learning the context specific causal structure using Algorithm 3 and 7

In order to assess the influence of the CI relations from the DAG in the final CStree, we tested the use of Algorithm 3 with the following experiment. First we generate a dataset from the above procedure, then we run Algorithm 3 on this dataset with and without using the CI relations from the CPDAG learnt in the first step. Below is an example where the DAG picks up a conditional independence relation which is not picked up by strictly running only the context specific independence tests.

Figure 4.5: CStree learnt using the CI relations from the DAG, which encodes the relation $X_3 \perp\!\!\!\perp X_2, X_1$

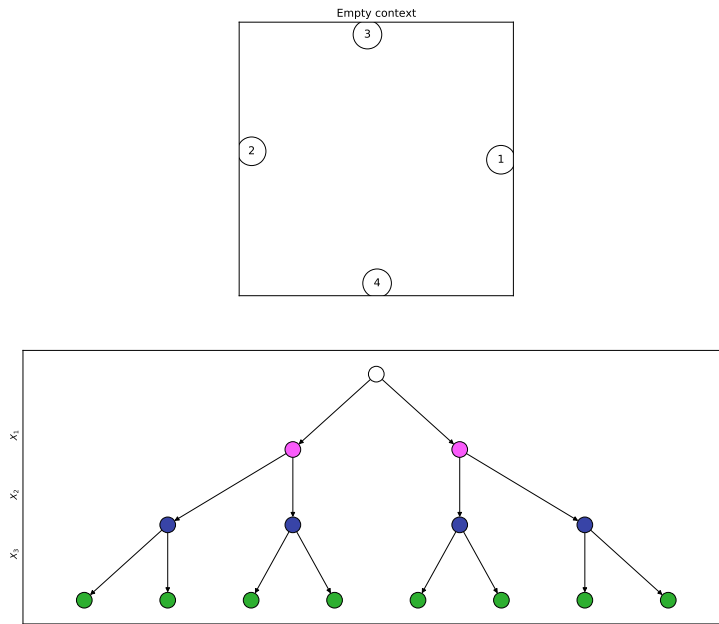
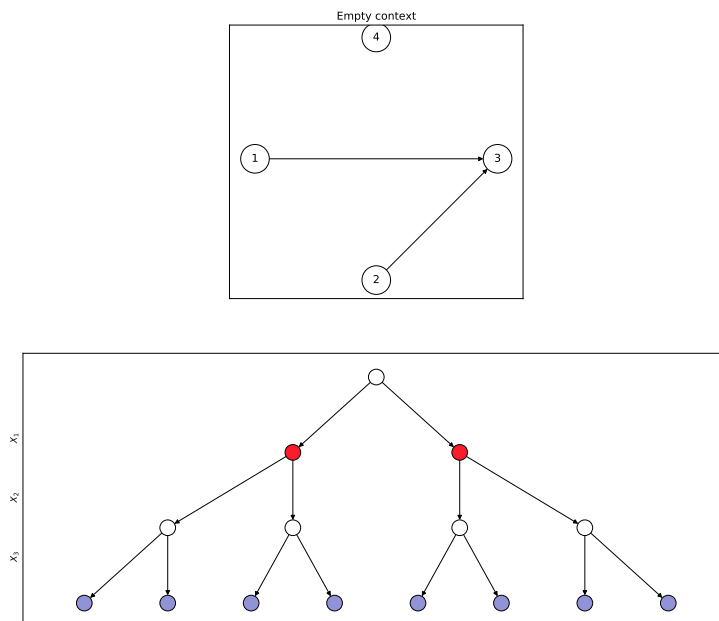


Figure 4.6: CStree learnt without the CI relations from the DAG, which does not encode the relation $X_3 \perp\!\!\!\perp X_2, X_1$ unlike when using the DAG on the same dataset



4.4 Coronary disease data

This is a dataset consisting of features which might increase the risk for coronary thrombosis, and consists of 1841 samples from men ¹. These features are detailed below.

Variable name	Node	Outcomes
Smoking	1	{Yes, No}
Strenuous mental work	2	{Yes, No}
Strenuous physical work	3	{Yes, No}
Systolic blood pressure	4	$\{(-\infty, 140], (140, -\infty)\}$
Ratio of beta and alpha lipoproteins	5	$\{(-\infty, 3], (3, -\infty)\}$
Coronary heart disease family history	6	{Yes, No}

¹ Z. Reinis, J. Pokorny, V. Basika, J. Tis-
erova, K. Gorican, D. Horakova,
E. Stuchlikova, T. Havranek, and
F. Hrabovsky. Prognostic significance
of the risk profile in the prevention of
coronary heart disease. *Bratisl Lek Listy*,
76:137–150, 1981

This data set does not require further pre-processing since the out-
come space is already categorical, and there were no missing values.

We first perform the following experiment: Learn a DAG model from
both the PC algorithm and the Hill climbing search algorithm. Now
for each topological ordering of all the DAGs which are Markov
equivalent, we store the following CStrees:

1. The CStree generated from the DAG itself
2. The CStree learnt only from context specific information without
encoding the CI relations from the DAG
3. The CStree learnt from encoding the CI relations from the DAG
and then learning further context specific information

For each of the cases above we create/merge stages using the
Anderson-Darling test ², the Epps-Singleton test ³, and the sym-
metric KL divergence with different threshold values. The use of the
symmetric KL divergence is motivated by the unreliability of condi-
tional independence testing in practice. Given 2 distributions \mathbb{P}, \mathbb{Q}
over the same discrete state space \mathcal{X} , the symmetric KL divergence is
as follows.

$$D_{SKL}(\mathbb{P}, \mathbb{Q}) = \sum_{x \in \mathcal{X}} \mathbb{P}(X = x) \log \frac{\mathbb{P}(X = x)}{\mathbb{Q}(X = x)} + \mathbb{Q}(X = x) \log \frac{\mathbb{Q}(X = x)}{\mathbb{P}(X = x)}$$

² F. W. Scholz and M. A. Stephens. K-
sample anderson-darling tests. *Journal
of the American Statistical Association*,
82(399):918–924, 1987

³ T.W. Epps and Kenneth J. Singleton.
An omnibus test for the two-sample
problem using the empirical charac-
teristic function. *Journal of Statistical
Computation and Simulation*, 26(3-4):177–
203, 1986

We apply Algorithm 3 to the setup above, and record the CStrees
with the minimum number of stages over all possible causal order-
ings, alongside their BIC scores. If we have CStrees with the same
number of minimum stages, we record the highest BIC score among

Goal: Min stages		Anderson		Epps		SKL 5×10^{-5}		SKL 5×10^{-6}		SKL 5×10^{-7}	
		Stages	BIC	Stages	BIC	Stages	BIC	Stages	BIC	Stages	BIC
PC	DAG	18	-6739.4	18	-6739.4	18	-6739.4	18	-6739.4	18	-6739.4
	CSTree w DAG	6	-6793.7	13	-6812.1	23	-6757.1	31	-7032.3	44	-6753.4
	CSTree w/o DAG	8	-6771.3	11	-6780.4	15	-6775.4	18	-6739.4	18	-6739.4
HC	DAG	31	-6733.3	31	-6733.3	31	-6733.3	31	-6733.3	31	-6733.3
	CSTree w DAG	8	-6786.5	31	-6837.3	26	-6799.0	39	-6828.0	47	-6755.9
	CSTree w/o DAG	8	-6788.7	15	-6809.7	26	-6801.4	23	-6805.4	30	-6729.5

them. We call this **experimental setup 1**, and the results of this experiment for the coronary dataset is summarized in the Table 4.2.

Here we see that the minimum number of stages over all the configurations is 6, and which results in the CStree shown in Figure 4.7. The causal ordering is 346125, and we get this CStree by learning the original DAG using the PC algorithm, followed by the Anderson-Darling test to learn context specific information without encoding the CI relations from the original DAG. However it can be seen that the BIC score is lower compared to the CStree generated from the DAG itself and the CStree learnt using the CI relations in the DAG before learning the context specific information. The latter CStree also has more stages than the CStree learnt without using the CI relations, which can be explained by the possibility that the DAG encoded CI relation(s) may have put nodes in the same stage, which otherwise may have been put in the same stage from the context specific testing procedure. Among the trees with the lowest stages, the CStree with the highest BIC score is generated with the ordering 534126, using the symmetric KL divergence with a threshold of $5e^{-5}$ for the staging procedure, from the DAG learnt from the Hill Climbing algorithm and without encoding the CI relations from it. It shown in Figure 4.8 We see that this CStree has 27 stages, which motivates us to define the **experimental setup 2**, where we aim to find the CStree with the

Table 4.2: Experimental setup 1 on the coronary dataset. Each cell under the column "Stages" is the minimum number of stages amongst all causal orderings generated from the either the PC or Hill Climbing algorithm, for different staging criteria. Many orderings produce CStrees with the same number of stages, in which case the BIC score in the table above is the maximum over those CStrees.

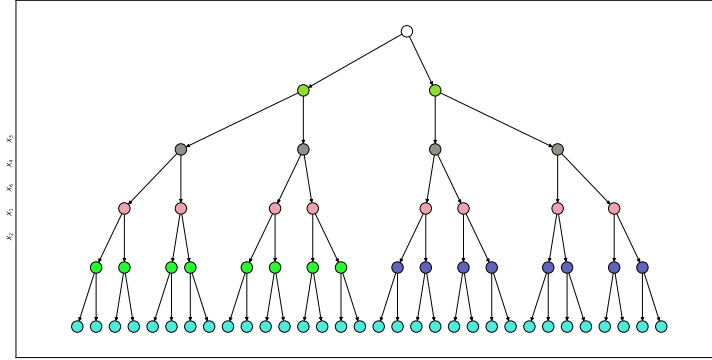


Figure 4.7: CStree for the coronary dataset with the lowest number of stages, which has ordering 346125.

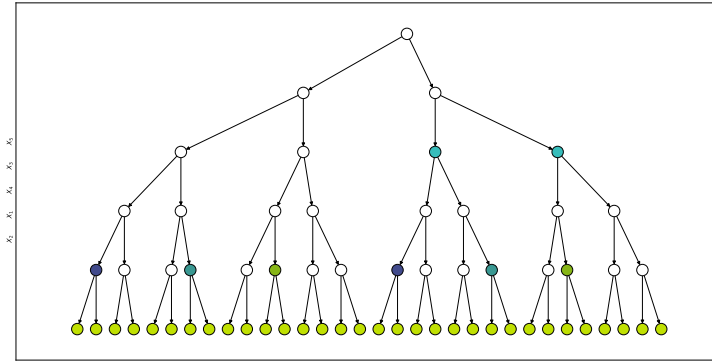


Figure 4.8: CStree with the highest BIC score among all CStrees with the minimum number of stages in experimental setup 1, which has ordering 534126.

highest BIC score instead of the minimum number of stages as in experimental setup 1. We show the results of this in Table 4.3.

We note that using the PC algorithm gave a MEC with 12 orderings, whilst the Hill climbing algorithm gave us a DAG whose MEC gave 2 orderings.

In order to compare our CStrees with similar staged trees, we run another experiment. First take the BIC optimal CStree that we have learnt which does not correspond to a DAG model. Take the ordering of this tree, and then fit corresponding staged trees using score based learning algorithms, optimizing for the BIC score. The algorithms used to fit the staged trees are the following:

1. Hill climbing - Start from the full independence model, then in a greedy manner, for each node, either assign it an existing stage or create a new stage depending on whichever option best increases the score.

Goal: Min stages		Anderson		Epps		SKL 5×10^{-5}		SKL 5×10^{-6}		SKL 5×10^{-7}	
		Stages	BIC	Stages	BIC	Stages	BIC	Stages	BIC	Stages	BIC
PC	DAG	18	-6739.4	18	-6739.4	18	-6739.4	18	-6739.4	18	-6739.4
	CSTree w DAG	7	-6783.7	13	-6812.1	23	-6757.1	43	-6749.6	44	-6753.4
	CSTree w/o DAG	8	-6771.3	11	-6780.4	15	-6775.4	18	-6739.4	18	-6739.4
HC	DAG	31	-6733.3	31	-6733.3	31	-6733.3	31	-6733.3	31	-6733.3
	CSTree w DAG	8	-6786.5	31	-6837.3	26	-6799.0	39	-6828.0	41	-6755.9
	CSTree w/o DAG	8	-6788.7	15	-6809.7	26	-6801.4	23	-6805.4	30	-6729.5

Table 4.3: Experimental setup 2 on the coronary dataset, the only difference to experimental setup 1 is that we now maximize the BIC score over all possible causal orderings.

2. Backward hill climbing - Similar to hill climbing except we start from a full model, i.e. one where each node belongs to a singleton stage.
3. Backward joining - Start from the full model, and iteratively merge stages depending on whether a certain measure between them is below a threshold. We use the KL divergence with twice the threshold which gave us the BIC optimal CStree from our methods.
4. Hierarchical clustering - Start from a full model, and iteratively cluster nodes in level k to c_k clusters using hierarchical clustering. We use the total variation measure alongside 2 clusters for each level. We call our approach to fit the CStree modified Algorithm 3 since the BIC optimal CStree is not necessarily learnt using the PC algorithm to get the initial DAG.

We call this **experimental setup 3**, and detail the results of it on the coronary dataset below. None of the staged trees are CStreets, and they are included in the Appendix.

Model	Algorithm	BIC	Stages
CStree	Modified Algorithm 3	-6729.5	30
Staged tree	Hill climbing	-6645.3	11
Staged tree	Backward hill climbing	-6641.1	13
Staged tree	Backward joining	-6790.4	57
Staged tree	Hierarchial clustering	-6654.7	10

In order to generate the sequence of DAGs that capture the context specific causal information from this dataset, we require 4 variables for computational feasibility. We do this by first removing the variable that gives us the the best score after its removal, which in this case is the variable corresponding to systolic blood pressure. We then generate the 5 BIC optimal CStrees arising after removing each of the remaining 5 variables.

4.5 Mice protein expression data

This is a dataset with expression levels of 77 proteins, measured in the cerebral cortex of 8 classes of mice ⁴. There are 38 control mice and 34 trisomic mice, and for each of them 15 separate measurements were taken which gives a total of 1080 samples. The aim is to identify features that can discriminate between the 8 classes that make up all possible combinations of the following 3 binary features - Control or Trisomy, stimulated to learn or not, injected with saline or memantine. This dataset had missing values, which is why we handle by first inspecting the missing value counts for each feature and then removing features which had 180 or more missing values. This removes the expression data corresponding to the columns BAD_N , $BCL2_N$, $H3AcK18_N$, $EGR1_N$, $H3MeK4_N$. This still leaves 72 features, which we further reduce by recursive feature elimination ⁵ where we first train a linear support vector classifier on all the features and prune the least important features in a greedy manner until we arrive at the number of features we want. After this we compute the medians for each feature and assign each feature to take binary outcomes depending on whether or not the value is above or below the median.

We choose 7 features and the predictor variable, giving a system of 8 variables. This number is chosen since this gives a maximum of 196 causal orderings which is computationally feasible, meanwhile increasing the variables to 9 and 10 increase this to 1008 and 13608 respectively. These are detailed below.

Variable name	Node	Outcomes
AKT_N	1	$\{(-\infty, Q_2^1], (Q_2^1, \infty)\}$
APP_N	2	$\{(-\infty, Q_2^2], (Q_2^2, \infty)\}$
$SOD1_N$	3	$\{(-\infty, Q_2^3], (Q_2^3, \infty)\}$
$NR2B_N$	4	$\{(-\infty, Q_2^4], (Q_2^4, \infty)\}$
$pNUMB_N$	5	$\{(-\infty, Q_2^5], (Q_2^5, \infty)\}$
$IL1B_N$	6	$\{(-\infty, Q_2^6], (Q_2^6, \infty)\}$
SYP_N	7	$\{(-\infty, Q_2^7], (Q_2^7, \infty)\}$
Class	8	$\{C_i\}_{i=1}^8$

Here, Q_2^i refers to the median of the variable corresponding to node i .

The results from Table 4.4 show once again that the Anderson-Darling test tends to merge stages more often thus resulting in a CStree with the minimum number of stages at the expense of a worse fit to the data, as indicated by the BIC score. We show the result of experimental setup 2 on this dataset in Table 4.5

⁴ Clara Higuera, Katheleen J. Gardiner, and Krzysztof J. Cios. Self-organizing feature maps identify proteins critical to learning in a mouse model of down syndrome. *PLOS ONE*, 10(6):e0129126, June 2015

⁵ Isabelle Guyon, Jason Weston, Stephen Barnhill, and Vladimir Vapnik. *Machine Learning*, 46(1/3):389–422, 2002

	Goal: Min Stages	Anderson		Epps		SKL 5×10^{-4}		SKL 5×10^{-5}		SKL 5×10^{-6}	
		Stages	BIC	Stages	BIC	Stages	BIC	Stages	BIC	Stages	BIC
PC	DAG	154	-5615.9	154	-5615.9	154	-5615.9	154	-5615.9	154	-5615.9
	CStree w DAG	7	-6555.4	8	-6442.3	61	-5678.9	64	-5640.6	78	-5587.6
	CStree w/o DAG	11	-6591.4	11	-6591.4	67	-5702.0	68	-5705.4	71	-5692.2
HC	DAG	102	-5376.6	102	-5376.6	102	-5376.6	102	-5376.6	102	-5376.6
	CStree w DAG	13	-6373.8	82	-6095.1	33	-6017.8	41	-5857.2	41	-5857.2
	CStree w/o DAG	12	-6404.1	20	-6179.9	41	-5852.3	49	-5815.0	59	-5837.1

Table 4.4: Experimental setup 1 on the mice protein expression dataset

	Goal: Max BIC	Anderson		Epps		SKL 5×10^{-4}		SKL 5×10^{-5}		SKL 5×10^{-6}	
		Stages	BIC	Stages	BIC	Stages	BIC	Stages	BIC	Stages	BIC
PC	DAG	154	-5615.9	154	-5615.9	154	-5615.9	154	-5615.9	154	-5615.9
	CSTree w DAG	15	-6542.4	45	-5804.4	113	-5474.4	115	-5513.2	133	-5433.0
	CSTree w/o DAG	11	-6591.4	36	-5906.2	75	-5598.9	76	-5602.4	79	-5589.1
HC	DAG	102	-5376.6	102	-5376.6	102	-5376.6	102	-5376.6	102	-5376.6
	CSTree w DAG	13	-6373.8	85	-6019.9	41	-5857.2	41	-5857.2	41	-5857.2
	CSTree w/o DAG	13	-6387.2	21	-6149.2	41	-5852.3	49	-5815.0	59	-5837.1

Table 4.5: Experimental setup 2 on the mice protein expression dataset.

Model	Algorithm	BIC	Stages
CStree	Modified Algorithm 3	-5433.0	133
Staged tree	Hill climbing	-5320.1	28
Staged tree	Backward hill climbing	-5227.6	36
Staged tree	Backward joining	-6427.1	148
Staged tree	Hierarchial clustering	-5990.4	17

4.6 Supersymmetry data

This is a simulated dataset involving 8 kinematic features measured by detectors in a particle accelerator alongside 10 more high-level features which are functions of the original 8 features ⁶. The aim is to identify the class label for each sample which denotes whether the measurement corresponds to a signal or background event. The dataset contains in total 5000000 samples. There are no missing values in this dataset, and we apply it on the 8 low-level features, and after taking a random subsample of size 2500000 samples. Our aim with using this dataset is to assess the empirical feasibility of our methods on large sample sizes. The number of valid causal orderings varies largely depending on the subsampled data, which ranged from 200 to 8000 in our experiments. As a result, we simply learn the CStree with the first ordering in the algorithm, using the KL divergence with a threshold of 5×10^{-7} . The process of learning the DAG model, converting it to a CStree, and learning CStrees using with and without the CI relations from the DAG for one ordering took approximately 6-7 hours on an Intel i7-7700 machine.

⁶ P. Baldi, P. Sadowski, and D. Whiteson. Searching for exotic particles in high-energy physics with deep learning. *Nature Communications*, 5(1), July 2014

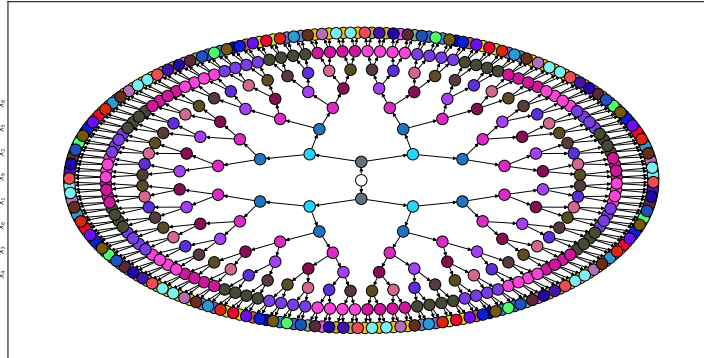


Figure 4.9: First CStree learnt from the supersymmetry dataset, with causal ordering 852916347.

The resulting CStree does not contain any non-singleton stage, which is perhaps surprising considering that one might expect many causal interactions to be present, and thus more stages.

4.7 Vitamin-D dataset

This dataset comes from a real study on Vitamin D and mortality ⁷, and contains 4 feature variables corresponding to age, a binary indicator denoting mutations in the filaggrin gene, vitamin D levels measured as serum 25-OH-D (nmol/L), follow up time, alongside the predictor variable which is a binary value indicating knowledge

⁷ Torben Martinussen, Ditte Nørbo Sørensen, and Stijn Vansteelandt. Instrumental variables estimation under a structural cox model. *Biostatistics*, 20(1):65–79, November 2017

Goal: Min stages		Anderson		Epps		SKL 5×10^{-1}		SKL 5×10^{-2}		SKL 5×10^{-4}	
		Stages	BIC	Stages	BIC	Stages	BIC	Stages	BIC	Stages	BIC
PC	DAG	12	-8410.1	12	-8410.1	12	-8410.1	12	-8410.1	12	-8410.1
	CSTree w DAG	4	-9041.5	13	-8800.0	4	-9041.5	5	-8771.7	9	-8532.3
	CSTree w/o DAG	6	-8480.1	12	-8410.1	4	-9041.5	9	-8405.2	12	-8410.1
HC	DAG	26	-8462.3	26	-8462.3	26	-8462.3	26	-8462.3	26	-8462.3
	CSTree w DAG	4	-9041.5	36	-8871.7	4	-9041.5	5	-8771.7	9	-8526.3
	CSTree w/o DAG	5	-8771.7	11	-8701.7	4	-9041.5	11	-8413.0	14	-8418.0

Table 4.6: Experimental setup 1 on the Vitamin D dataset

of whether the subject passed away during the follow up. Looking at the data it is clear that the ages are grouped into 4 distinct groups. We group the vitamin D measurement into 4 groups based on the quartiles of the data, and the follow up time is grouped into a binary outcome with the median being the cutoff point. The resulting dataset is summarized in the table below.

Variable name	Node	Outcomes
Age	1	$\{[40, 47), [47, 57), [57, 67), [67, 80)\}$
Filaggrin	2	$\{Yes, No\}$
Vitamin D	3	$\{[-\infty, Q_1^3), [Q_1^3, Q_2^3), [Q_2^3, Q_3^3), [Q_3^3, -\infty)\}$
Follow up time	4	$\{(-\infty, Q_2^4], (Q_2^4, \infty)\}$
Passed away	5	$\{Yes, No\}$

Here Q_1^i, Q_2^i, Q_3^i respectively denote the lower quartile, median and upper quartile of the variable corresponding to node i .

We perform the experimental setup 1 and 2 on this dataset and get the following results.

Table 4.8 shows the results of experimental setup 3 on the Vitamin D dataset.

Here we see that the CStree outperforms all staged tree models learnt in terms of the BIC score, and also provides the sparsest model if

Goal: Max BIC		Anderson		Epps		SKL 5×10^{-1}		SKL 5×10^{-2}		SKL 5×10^{-4}	
		Stages	BIC	Stages	BIC	Stages	BIC	Stages	BIC	Stages	BIC
PC	DAG	12	-8410.1	12	-8410.1	12	-8410.1	12	-8410.1	12	-8410.1
	CSTree w DAG	6	-8487.9	14	-8500.0	5	-8749.9	19	-8428.2	37	-8489.0
	CSTree w/o DAG	6	-8480.1	12	-8410.1	5	-8749.9	9	-8405.2	12	-8410.1
HC	DAG	26	-8462.3	26	-8462.3	26	-8462.3	26	-8462.3	26	-8462.3
	CSTree w DAG	5	-8779.6	36	-8871.7	4	-9041.5	8	-8487.9	13	-8525.8
	CSTree w/o DAG	5	-8771.7	11	-8701.7	4	-9041.5	11	-8413.0	14	-8418.0

Table 4.7: Experimental setup 2 on the Vitamin D dataset

Model	Algorithm	BIC	Stages
CStree	Modified Algorithm 3	-8405.2	9
Staged tree	Hill climbing	-8423.1	11
Staged tree	Backward hill climbing	-8423.4	12
Staged tree	Backward joining	-8560.8	26
Staged tree	Hierarchial clustering	-8910.0	10

Table 4.8: Experimental setup 3 on Vitamin D dataset with ordering 53421

we exclude the result of the hierarchical clustering algorithm where we pre define the number of maximum stages per level. However, it is very important to note that the causal ordering is starting from the variable corresponding to whether the subject passed away or not, which is not consistent with the true causal ordering. This is because the outcome corresponding to this variable is the last to be measured in the data generating process. This suggests errors in the initial DAG, which may have been caused due to instability in the conditional independence tests or being stuck at a local minima. Similarly, it would make sense if the variable corresponding to the follow up time happens second to last, since the only measurement taken during the follow up is whether the subject passed away or not. We thus learn the CStrees for all possible orderings ending with 45, and rerun experimental setup 3 with a more realistic causal ordering. The results are summarized in Table 4.9 below.

Model	Algorithm	BIC	Stages
CStree	Modified Algorithm 3	-8776.5	9
Staged tree	Hill climbing	-8406.5	13
Staged tree	Backward hill climbing	-8403.5	16
Staged tree	Backward joining	-8547.1	19
Staged tree	Hierarchical clustering	-8628.6	10

Table 4.9: Experimental setup 3 on Vitamind D dataset with ordering 32145

As anticipated, none of the DAGs in the MEC of the DAG learnt using both the PC algorithm and the Hill climbing search were consistent with any of the orderings ending with 45. Thus, we only rely on context specific information learnt from the staging procedure, without any CI relations. Technically, the highest BIC score was achieved using the Epps test, which gives the CStree with ordering 12345 a BIC score of -8734.5 however this CStree has 48 stages. This is followed by the CStree obtained using the symmetric KL divergence with a threshold of 5×10^{-2} , which gives the ordering 21345 with a BIC score of -8771.7. However, this is a DAG model with one edge $4 \rightarrow 5$ in the empty context graph. We found changing this threshold to 5×10^{-3} gave us a slightly worse BIC score but a CStree that encodes interesting context specific information, which is why we choose this ordering as a basis for comparison with staged trees. We show this CStree alongside its minimal context DAGs below.

4.8 Discussion of empirical performance

We see in general that using the symmetric KL divergence to identify stages in the CStree result a higher BIC score when compared to choosing a CStree with the minimum number of stages. There are many reasons which could lead to this happening in practice, with

the first being unreliability of the statistical independence tests, which is in general a very hard task ⁸.

If the data generating distribution happens to be faithful to a CStree with a large number of stages, particularly non-singleton stages, then it is possible that there is not enough samples to learn the context specific information in the finite data setting. We note that a lot of this information was learnt by comparing samples with extremely different sizes, so in practice it might be a good idea to merge the corresponding nodes into the same stage if a certain criteria is met, for example, only if both sample sizes are at least half of their mean.

The instability of DAG learning algorithms used to get the causal orderings also play a crucial role, and we observed this in practice. For example, we observed that if we happen to know the causal ordering and then feed it to the algorithm, we could learn a DAG in the initial step whose MEC contains no DAGs which are consistent with this causal ordering. This may also present itself to be a problem if we have partial knowledge of the ordering, for example in the mice cortex data experiment, it is plausible to think of the predictor class to be a function of the gene expression levels.

One important aspect of this process is to choose the best CStree from all possible causal ordering, which is a model selection problem. In Algorithm 3 we opt to choose the one with the fewest stages since it relates to choosing the model with the fewest parameters, i.e. the simplest model. This goes in hand with the principle of Occam's razor - in the face of many possible models, choose the simplest model. Here a simple model refers to one which has few parameters. In practice however we see that problems arising due to misleading conditional independence tests and small data samples might lead to the simplest model being too simple, compromising the fit to the data.

⁸ Rajen D. Shah and Jonas Peters. The hardness of conditional independence testing and the generalised covariance measure. *The Annals of Statistics*, 48(3), Jun 2020

Conclusions

5.1 Summary

We start with DAGs as a means to encode CI relations, and how one can use the characterization of Markov Equivalence in DAGs to learn causal structure through CI testing. We then cover the limitations of CI relations in comparison to CSI relations, and go over CStrees as a means to encode such CSI relations. We then show how one can learn these CSI relations from observational data and learn a CStree, and how to compute minimal contexts which are important when it comes to visualizing higher dimensional CStrees. We then apply these techniques to synthetic and real data.

5.2 Future work

One of the more natural extensions of this work is to learn CStrees from interventional data. This can already be done with DAGs ¹. On the topic of model selection, it might be interesting to see the applicability of Bayesian model selection for CStrees, whereby the model evidence is used as the basis for comparing models which automatically penalizes over-complex models whilst also penalizing models which do not agree with the observed data ². Similar approaches have been applied to Gaussian DAG models ³. Generalization to missing data problems would also be an interesting avenue, considering that the PC algorithm has been recently extended for missing data instances in DAGs ⁴. But perhaps the most impactful extension would be formulating the problem of finding the best CStree into a continuous optimization problem, which can lead to scalable score based methods. Recent work has formulated the problem of searching for the best DAG according to some metric by using a characterization of acyclicity that is smooth and exact, allowing the conversion of

¹ Karren Yang, Abigail Katcoff, and Caroline Uhler. Characterizing and learning equivalence classes of causal DAGs under interventions. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5541–5550. PMLR, 10–15 Jul 2018

² David J. C. MacKay. *Bayesian Interpolation*, pages 39–66. Springer Netherlands, Dordrecht, 1992

³ Federico Castelletti. Bayesian model selection of gaussian directed acyclic graph structures. *International Statistical Review*, 88(3):752–775, June 2020

⁴ Ruibo Tu, Cheng Zhang, Paul Ackermann, Karthika Mohan, Hedvig Kjellström, and Kun Zhang. Causal discovery in the presence of missing data. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pages 1762–1770. PMLR, 16–18 Apr 2019

the combinatorial problem into a purely continuous problem ⁵.

⁵ Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. Dags with no tears: Continuous optimization for structure learning. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018

Bibliography

- [1] Biwei Huang, Kun Zhang, Mingming Gong, and Clark Glymour. Causal discovery and forecasting in nonstationary environments with state-space models, 2019.
- [2] Pengfei Hu, Rong Jiao, Li Jin, and Momiao Xiong. Application of causal inference to genomic analysis: Advances in methodology. *Frontiers in Genetics*, 9:238, 2018.
- [3] Jakob Runge, Sebastian Bathiany, Erik Bollt, Gustau Camps-Valls, Dim Coumou, Ethan Deyle, Clark Glymour, Marlene Kretschmer, Miguel D Mahecha, Jordi Muñoz-Marí, et al. Inferring causation from time series in earth system sciences. *Nature communications*, 10(1):1–13, 2019.
- [4] Judea Pearl and Dana Mackenzie. *The Book of Why: The New Science of Cause and Effect*. Basic Books, Inc., USA, 1st edition, 2018.
- [5] Duligur Ibeling Thomas Icard Elias Bareinboim, Juan D. Correa. On pearl’s hierarchy and the foundations of causal inference, 2020.
- [6] R.A. Fisher. *The design of experiments*. 1935. Oliver and Boyd, Edinburgh, 1935.
- [7] David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of machine learning research*, 3(Nov):507–554, 2002.
- [8] Peter Spirtes and Clark Glymour. An algorithm for fast recovery of sparse causal graphs. *Social Science Computer Review*, 9(1):62–72, 1991.

- [9] L Solus, Y Wang, and C Uhler. Consistency Guarantees for Greedy Permutation-Based Causal Inference Algorithms. *Biometrika*, 01 2021. asaa104.
- [10] Ioannis Tsamardinos, Laura E Brown, and Constantin F Aliferis. The max-min hill-climbing bayesian network structure learning algorithm. *Machine learning*, 65(1):31–78, 2006.
- [11] Rodrigo A Collazo, Christiane Görgen, and Jim Q Smith. *Chain event graphs*. CRC Press, 2018.
- [12] Tomi Silander and Tze-Yun Leong. A dynamic programming algorithm for learning chain event graphs. In *International Conference on Discovery Science*, pages 201–216. Springer, 2013.
- [13] Peter Thwaites, Jim Q. Smith, and Eva Riccomagno. Causal analysis with chain event graphs. *Artificial Intelligence*, 174(12):889–909, 2010.
- [14] Christiane Görgen and Jim Q. Smith. Equivalence classes of staged trees, 2017.
- [15] Bernhard Schölkopf. *Causality for machine learning*, 2019.
- [16] Edward H Simpson. The interpretation of interaction in contingency tables. *Journal of the Royal Statistical Society: Series B (Methodological)*, 13(2):238–241, 1951.
- [17] Clifford H. Wagner. Simpson’s paradox in real life. *The American Statistician*, 36(1):46–48, 1982.
- [18] David H. Wolpert. What is important about the no free lunch theorems?, 2020.
- [19] Clark Glymour, Kun Zhang, and Peter Spirtes. Review of causal discovery methods based on graphical models. *Frontiers in Genetics*, 10(nil):nil, 2019.
- [20] Eliana Duarte and Liam Solus. Algebraic geometry of discrete interventional models, 2020.
- [21] Steffen L Lauritzen. *Graphical models*, volume 17. Clarendon Press, 1996.
- [22] Christopher Meek. Strong completeness and faithfulness in bayesian networks. *arXiv preprint arXiv:1302.4973*, 2013.
- [23] Tom S. Verma and Judea Pearl. On the equivalence of causal models. *CoRR*, abs/1304.1108, 2013.

- [24] Caroline Uhler, Garvesh Raskutti, Peter Bühlmann, and Bin Yu. Geometry of the Faithfulness Assumption in Causal inference. *The Annals of Statistics*, 41(2):436 – 463, 2013.
- [25] P. Spirtes, C. Glymour, and R. Scheines. *Causation, Prediction, and Search*. MIT press, 2nd edition, 2000.
- [26] Markus Kalisch and Peter Bühlmann. Estimating high-dimensional directed acyclic graphs with the pc-algorithm. *Journal of Machine Learning Research*, 8(22):613–636, 2007.
- [27] Christopher Meek. Causal inference and causal explanation with background knowledge, 2013.
- [28] Steen A. Andersson, David Madigan, and Michael D. Perlman. A characterization of markov equivalence classes for acyclic digraphs. *The Annals of Statistics*, 25(2), April 1997.
- [29] Eliana Duarte and Liam Solus. Representation of context-specific causal models with observational and interventional data, 2021.
- [30] Robert Endre Tarjan. Edge-disjoint spanning trees and depth-first search. *Acta Informatica*, 6(2):171–185, 1976.
- [31] Kayvan Sadeghi and Steffen Lauritzen. Markov Properties for Mixed graphs. *Bernoulli*, 20(2):676 – 696, 2014.
- [32] Thomas Verma and Judea Pearl. *Causal Networks: Semantics and Expressiveness*. "Graphoids: A Computer Representation for Dependencies and Relevance in Automated Reasoning (Computer Information Science).", pages 69–76. Uncertainty in Artificial Intelligence. Elsevier, 1990.
- [33] Judea Pearl. *Causality*. Cambridge University Press, 2009.
- [34] Z. Reinis, J. Pokorny, V. Basika, J. Tiserova, K. Gorican, D. Horakova, E. Stuchlikova, T. Havranek, and F. Hrabovsky. Prognostic significance of the risk profile in the prevention of coronary heart disease. *Bratisl Lek Listy*, 76:137–150, 1981.
- [35] F. W. Scholz and M. A. Stephens. K-sample anderson-darling tests. *Journal of the American Statistical Association*, 82(399):918–924, 1987.
- [36] T.W. Epps and Kenneth J. Singleton. An omnibus test for the two-sample problem using the empirical characteristic function. *Journal of Statistical Computation and Simulation*, 26(3-4):177–203, 1986.

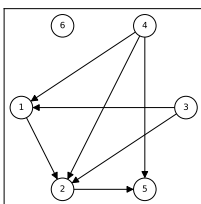
- [37] Clara Higuera, Katheleen J. Gardiner, and Krzysztof J. Cios. Self-organizing feature maps identify proteins critical to learning in a mouse model of down syndrome. *PLOS ONE*, 10(6):e0129126, June 2015.
- [38] Isabelle Guyon, Jason Weston, Stephen Barnhill, and Vladimir Vapnik. *Machine Learning*, 46(1/3):389–422, 2002.
- [39] P. Baldi, P. Sadowski, and D. Whiteson. Searching for exotic particles in high-energy physics with deep learning. *Nature Communications*, 5(1), July 2014.
- [40] Torben Martinussen, Ditte Nørbo Sørensen, and Stijn Vansteelandt. Instrumental variables estimation under a structural cox model. *Biostatistics*, 20(1):65–79, November 2017.
- [41] Rajen D. Shah and Jonas Peters. The hardness of conditional independence testing and the generalised covariance measure. *The Annals of Statistics*, 48(3), Jun 2020.
- [42] Karren Yang, Abigail Katcoff, and Caroline Uhler. Characterizing and learning equivalence classes of causal DAGs under interventions. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5541–5550. PMLR, 10–15 Jul 2018.
- [43] David J. C. MacKay. *Bayesian Interpolation*, pages 39–66. Springer Netherlands, Dordrecht, 1992.
- [44] Federico Castelletti. Bayesian model selection of gaussian directed acyclic graph structures. *International Statistical Review*, 88(3):752–775, June 2020.
- [45] Ruibo Tu, Cheng Zhang, Paul Ackermann, Karthika Mohan, Hedvig Kjellström, and Kun Zhang. Causal discovery in the presence of missing data. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pages 1762–1770. PMLR, 16–18 Apr 2019.
- [46] Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. Dags with no tears: Continuous optimization for structure learning. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

Appendix A : Implementation remarks

Implementation of this work was done in Python. Draft code repository is available at github.com/mnazaal/masters-thesis. This document was generated with org-mode. Use of scikit, pandas, numpy. Test driven development with pytest.

Appendix B : Extra figures

PC



Hill climbing

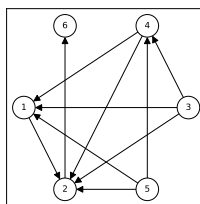
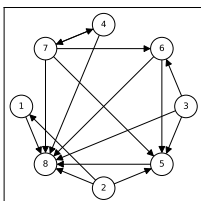


Figure 8.1: CPDAG representing the MEC of the DAG learnt from the coronary dataset.

PC



Hill climbing

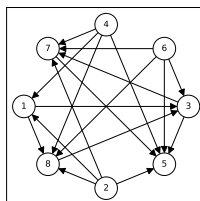
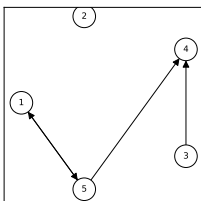


Figure 8.2: CPDAG representing the MEC of the DAG learnt from the mice protein expression dataset.

PC



Hill climbing

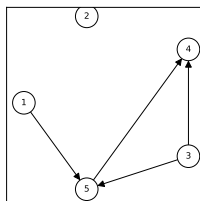
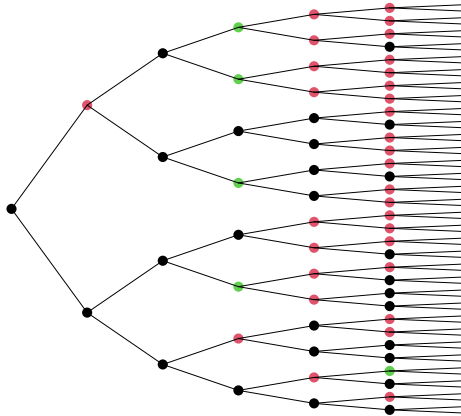
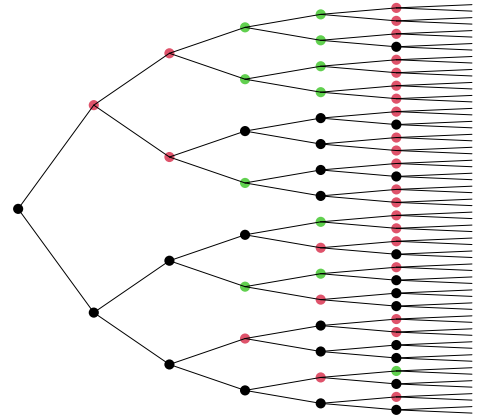


Figure 8.3: CPDAG representing the MEC of the DAG learnt from the Vitamin D dataset.

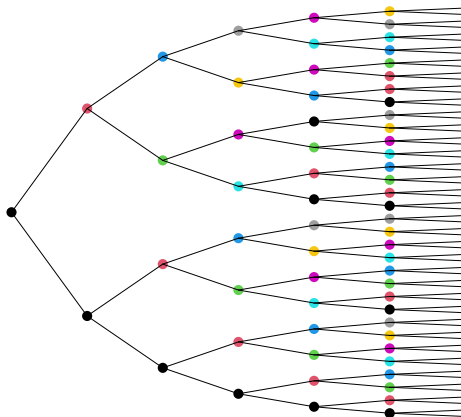
Hill Climbing



Backwards hill climbing



Backwards joining



Hierarchical clustering

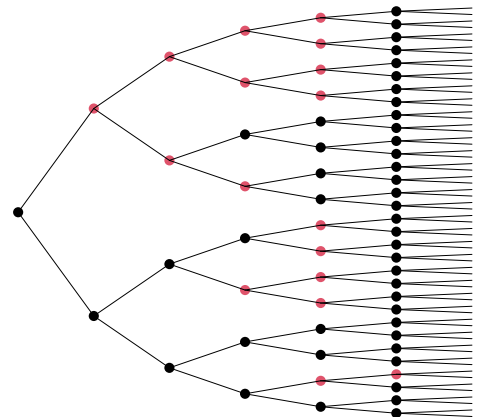
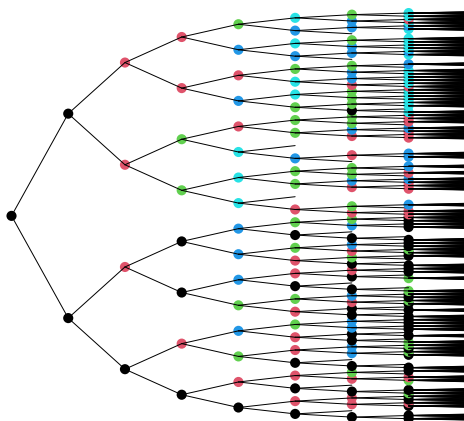
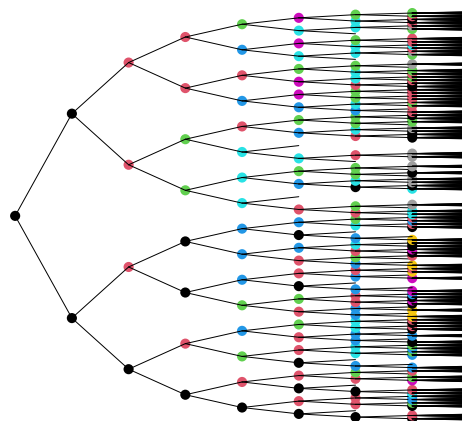


Figure 8.4: BIC Optimal staged trees learnt from the coronary dataset using the ordering 534126.

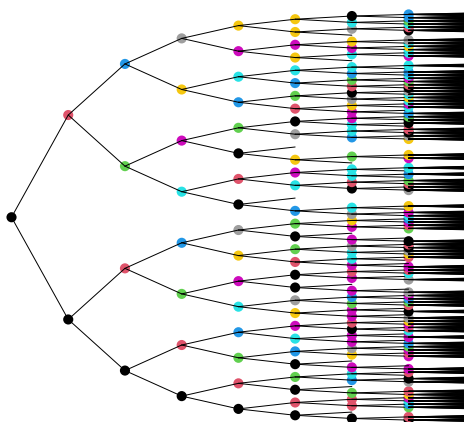
Hill Climbing



Backwards hill climbing



Backwards joining



Hierarchical clustering

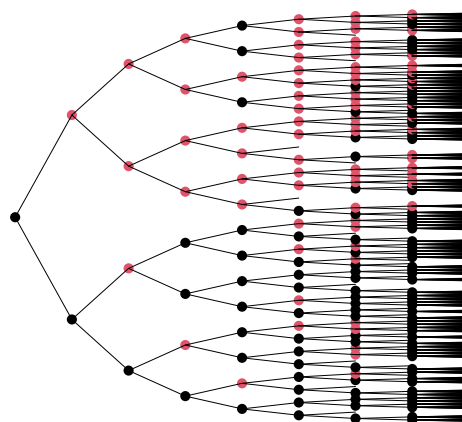
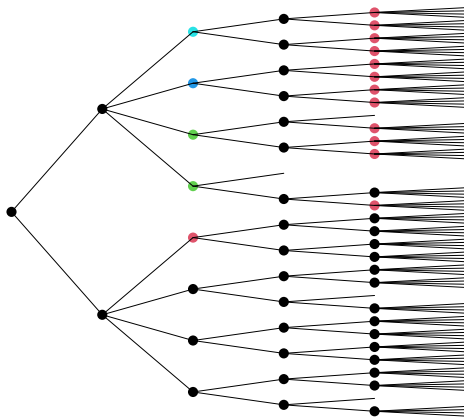
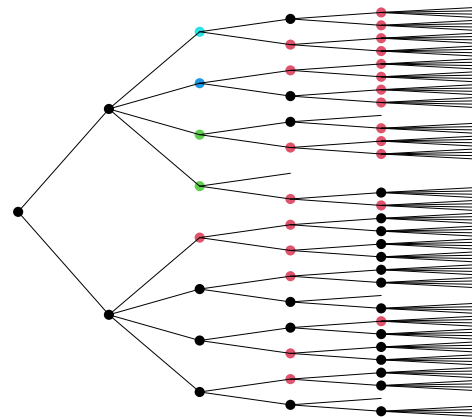


Figure 8.5: BIC Optimal staged trees learnt from the mice protein expression dataset using the ordering 32765148.

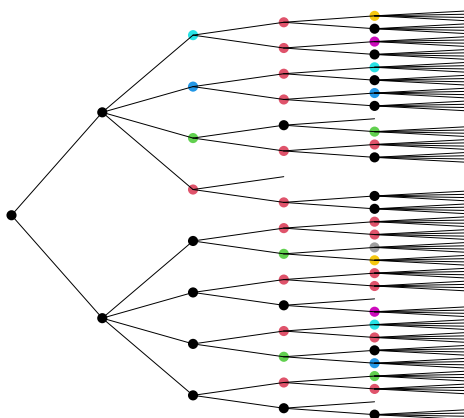
Hill Climbing



Backwards hill climbing



Backwards joining



Hierarchical clustering

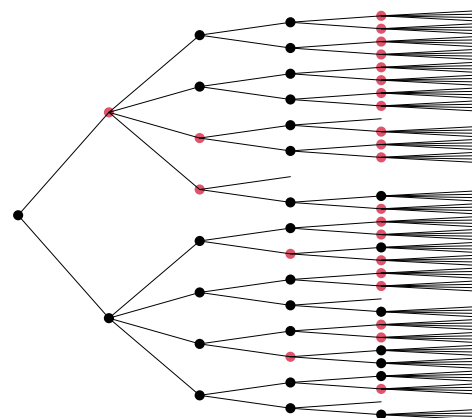
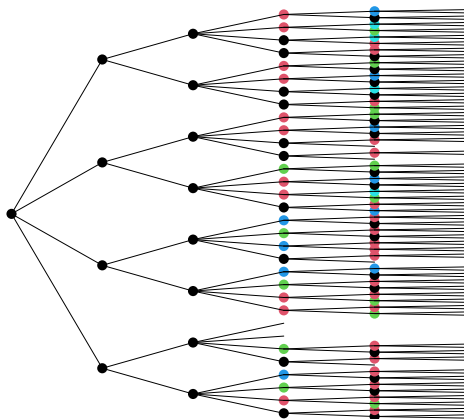
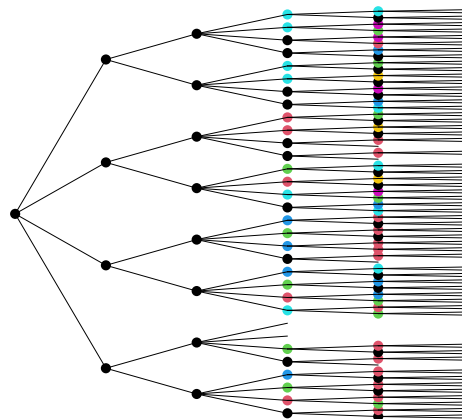


Figure 8.6: BIC Optimal staged trees learnt from the Vitamin D dataset using the ordering 53421.

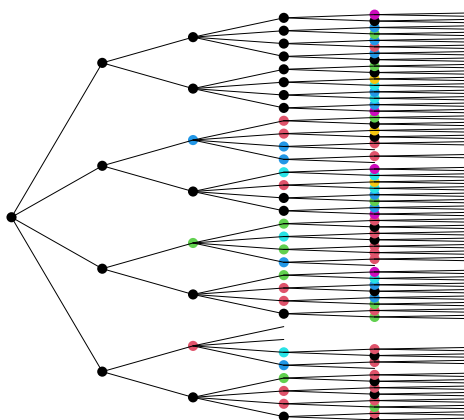
Hill Climbing



Backwards hill climbing



Backwards joining



Hierarchical clustering

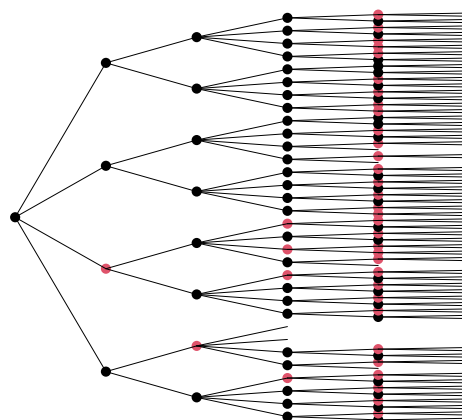


Figure 8.7: BIC Optimal staged trees learnt from the Vitamin D dataset using the ordering 32145.