

Abstract

Despite having a philosophical grounding from empiricism that spans 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 available data. 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 given Z and $C = a$ 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 tillgängliga data. 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 och $C = a$ 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>	7
1.1	<i>Motivation</i>	7
1.2	<i>Classical methods</i>	7
1.3	<i>Context specific causal discovery</i>	8
1.4	<i>Relevance to machine learning</i>	9
2	<i>Causal Discovery with Directed Acyclic Graphs</i>	10
2.1	<i>The Causal Discovery problem</i>	10
2.2	<i>Direct Acyclic Graphs (DAGs)</i>	11
2.3	<i>Causal Discovery Algorithms for DAGs</i>	14
2.4	<i>Limitations of using DAGs</i>	15
3	<i>Causal Discovery with Context Specific Trees</i>	16
3.1	<i>Context Specific Trees (CStrees)</i>	16
3.2	<i>Learning CStrees from observed data</i>	19
3.3	<i>Understanding high-dimensional CStrees</i>	23
3.4	<i>Model selection for CStrees</i>	36
4	<i>Experiments</i>	38
4.1	<i>Synthetic data</i>	38
4.2	<i>Coronary disease data</i>	42
4.3	<i>Mice protein expression data</i>	49
4.4	<i>Supersymmetry data</i>	53
4.5	<i>Vitamin-D data</i>	54
4.6	<i>Discussion of empirical performance</i>	58

5	<i>Conclusions</i>	61
5.1	<i>Summary</i>	61
5.2	<i>Future work</i>	61
6	<i>Bibliography</i>	63
7	<i>Appendix A : Implementation remarks</i>	68
8	<i>Appendix B : Extra figures</i>	69

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. Philosophical discussions of this topic date back to at least the Scottish enlightenment ¹, however the formal mathematical study and algorithmization of causality is a relatively recent enterprise. Whilst being an active research area for its own sake, causal discovery has many applications to diverse fields ranging from economics ² and genomics ³ to the climate sciences ⁴. 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 ⁵. In fact it has been recently shown that the strict distinction between the levels in this hierarchy hold in the measure theoretic sense ⁶.

1.2 Classical methods

Historically speaking the task of causal discovery was achieved by using clever experimental design ⁷. 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

¹ Hume, D. (1748). *Enquiry Concerning Human Understanding*

² Hünermund, P. and Bareinboim, E. (2019). Causal inference and data-fusion in econometrics. *arXiv preprint arXiv:1912.09104*

³ Hu, P., Jiao, R., Jin, L., and Xiong, M. (2018). Application of causal inference to genomic analysis: Advances in methodology. *Frontiers in Genetics*, 9:238

⁴ Runge, J., Bathiany, S., Bollt, E., Camps-Valls, G., Coumou, D., Deyle, E., Glymour, C., Kretschmer, M., Mahecha, M. D., Muñoz-Marí, J., et al. (2019). Inferring causation from time series in earth system sciences. *Nature communications*, 10(1):1–13

⁵ Pearl, J. and Mackenzie, D. (2018). *The Book of Why: The New Science of Cause and Effect*. Basic Books, Inc., USA, 1st edition

⁶ Elias Bareinboim, Juan D. Correa, D. I. T. I. (2020). On pearl's hierarchy and the foundations of causal inference

⁷ Fisher, R. (1935). *The design of experiments*. 1935. Oliver and Boyd, Edinburgh

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 whether smoking causes lung cancer which would involve forcing the subjects to smoke. This creates a need for more sophisticated methods involving the use of 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

Many of the causal models in wide use today make use of Directed Acyclic Graphs (DAGs), with the semantics that the directed edges represent direct 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 decades⁸, 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 information⁹, 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. 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 algorithm¹⁰, alongside applications to both synthetic and real world data.

⁸ Chickering, D. M. (2002b). Optimal structure identification with greedy search. *Journal of machine learning research*, 3(Nov):507–554; Spirtes, P. and Glymour, C. (1991). An algorithm for fast recovery of sparse causal graphs. *Social Science Computer Review*, 9(1):62–72; Solus, L., Wang, Y., and Uhler, C. (2021). Consistency Guarantees for Greedy Permutation-Based Causal Inference Algorithms. *Biometrika*. asaa104; and Tsamardinos, I., Brown, L. E., and Aliferis, C. F. (2006). The max-min hill-climbing bayesian network structure learning algorithm. *Machine learning*, 65(1):31–78

⁹ Collazo, R. A., Görgen, C., and Smith, J. Q. (2018). *Chain event graphs*. CRC Press; Silander, T. and Leong, T.-Y. (2013). A dynamic programming algorithm for learning chain event graphs. In *International Conference on Discovery Science*, pages 201–216. Springer; Thwaites, P., Smith, J. Q., and Riccomagno, E. (2010). Causal analysis with chain event graphs. *Artificial Intelligence*, 174(12):889–909; and

¹⁰ Spirtes, P. and Glymour, C. (1991). An algorithm for fast recovery of sparse causal graphs. *Social Science Computer Review*, 9(1):62–72

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, 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 under the assumption of independent and identically distributed (i.i.d.) data, and ignores information from interventions, domain shifts and temporal structure¹¹. As such, there are various problems which still require a causal model, which without it in some cases even give rise to seemingly nuanced paradoxes¹². One example is Simpson's paradox¹³, 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 well¹⁴.

¹¹ Schölkopf, B. (2019). Causality for machine learning. *arXiv preprint arXiv:1911.10500*

¹² Pearl, J. and Mackenzie, D. (2018). *The Book of Why: The New Science of Cause and Effect*. Basic Books, Inc., USA, 1st edition

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

¹⁴ Wagner, C. H. (1982). Simpson's paradox in real life. *The American Statistician*, 36(1):46–48

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 information 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,

¹

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

like the Bayesian Information 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 $G = (\mathbb{V}, \mathbb{E})$ which has no cycles.

Let $G = (\mathbb{V}, \mathbb{E})$ be a graph. We then have the following definitions. A node $u \in \mathbb{V}$ is a **parent** of another node $v \in \mathbb{V}$ if $(u, v) \in \mathbb{E}$, in this case we also say v is a **child** of u . If there is an edge $(u, v) \in \mathbb{E}$ or $(v, u) \in \mathbb{E}$ we say the nodes u, v are **adjacent**. A sequence of nodes (u_1, \dots, u_k) , $k \geq 2$ such that $(u_i, u_{i+1}) \in \mathbb{E}$ is called a path between u_1 and u_k . In this 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 $u \in \mathbb{V}$ 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.

Throughout this work we let $[p] = \{1, \dots, p\}$ be the set of nodes for DAGs.

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 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.

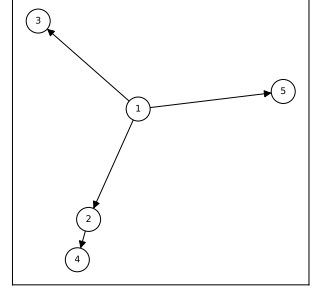


Figure 2.1: Example of a DAG $G = (\mathbb{V}, \mathbb{E})$ with $\mathbb{V} = \{1, 2, 3, 4, 5\}$ and $\mathbb{E} = \{(1, 2), (1, 3), (1, 5), (2, 4)\}$. Here $PA_G(2) = \{1\}$, $DS_G(2) = CH_G(2) = \{1\}$, $ND_G(2) = \{3, 5\}$

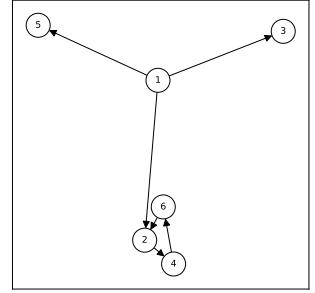


Figure 2.2: This graph is not a DAG since there is a cycle

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 (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 $(X_A \perp\!\!\!\perp_{\mathbb{P}} X_B \mid X_S)$) 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 $(X_A \perp\!\!\!\perp_{\mathbb{P}} X_B \mid X_S)$ 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.

Definition 2.4 (Local Markov property) Let \mathbb{G} be a DAG with nodes $[p]$. A probability distribution \mathbb{P} satisfies the local Markov property with respect to \mathbb{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, $(X_i \perp\!\!\!\perp_{\mathbb{P}} X_{ND_{\mathbb{G}}(i)} \mid X_{PA_{\mathbb{G}}(i)})$.

This formalizes the fact that in order to computationally generate data from a DAG \mathbb{G} , the value of each variable $X_i \in \mathcal{X}_i$ depends only on the values of the outcomes of its parents in \mathbb{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 $\sum_{i=1}^p |\mathcal{X}_i| \prod_{j \in PA_{\mathbb{G}}(i)} |\mathcal{X}_j|$ values. This is significantly smaller than having to store all possible probability values which would require $\sum_{i=1}^p |\mathcal{X}_i| \prod_{i=1}^p |\mathcal{X}_i|$ values. For binary variables assuming d parents for each variable, this is the difference between $p2^{d+1}$ and $p2^{p+1}$.

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

Definition 2.5 (Ordered Markov Property) Let \mathbb{G} be a DAG and $\pi = \pi_1 \cdots \pi_p$ a causal ordering of \mathbb{G} . A probability distribution \mathbb{P} satisfies the Ordered Markov property with respect to \mathbb{G} if we have $(X_i \perp\!\!\!\perp_{\mathbb{P}} X_{\{1, \dots, i-1\} \setminus PA_{\mathbb{G}}(i)} \mid X_{PA_{\mathbb{G}}(i)})$.

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

A distribution \mathbb{P} satisfying the local Markov property with respect to a DAG G is equivalent to that distribution satisfying the ordered 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 2.6 (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 2.7 (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 2.8 (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 states that the above notions are indeed equivalent⁵.

Theorem 2.9 (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 seen as structures that encode Conditional Independence (CI) relations. More importantly, d-separation encodes the

⁴ Recall that a path between 2 nodes is any sequence of edges connecting them irrespective of edge direction.

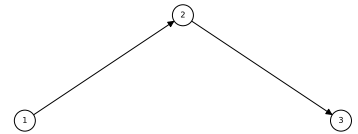


Figure 2.3: Chain

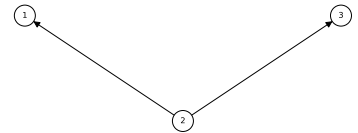


Figure 2.4: Fork/Common cause

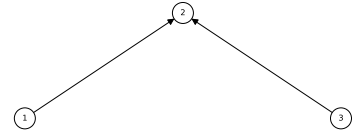


Figure 2.5: V-structure/ Collider/Immortality

⁵ Duarte, E. and Solus, L. (2020). Algebraic geometry of discrete interventional models. *arXiv preprint arXiv:2012.03593*

⁶ Lauritzen, S. L. (1996). *Graphical models*, volume 17. Clarendon Press

complete set of CI relations satisfied by all distributions Markov to a DAG, i.e. 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 2.10 (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 2.9 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 2.11 (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 2.10 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, then its v-structures, 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

⁷ Meek, C. (1995b). Strong completeness and faithfulness in bayesian networks. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence*, UAI'95, pages 411–418, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc; and Geiger, D., Verma, T., and Pearl, J. (1990). Identifying independence in bayesian networks. *Networks*, 20(5):507–534

⁸ Verma, T. S. and Pearl, J. (2013). On the equivalence of causal models. *CoRR*, abs/1304.1108

⁹ Meek, C. (1995b). Strong completeness and faithfulness in bayesian networks. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence*, UAI'95, pages 411–418, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc

¹⁰ Uhler, C., Raskutti, G., Bühlmann, P., and Yu, B. (2013). Geometry of the Faithfulness Assumption in Causal inference. *The Annals of Statistics*, 41(2):436 – 463

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

¹² Meek, C. (1995a). Causal inference and causal explanation with background knowledge. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence*, UAI'95, pages 403–410, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc

a is Complete PDAG (CPDAG) if every directed edge exists also in every DAG in the 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 graphs ¹³.

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 2.12 (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 conditionally independent given S , in the context $X_C = x_C$ (denoted $(X_A \perp\!\!\!\perp_{\mathbb{P}} X_B \mid X_S, X_C = x_C)$) 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 Andersson, S. A., Madigan, D., and Perlman, M. D. (1997). A characterization of markov equivalence classes for acyclic digraphs. *The Annals of Statistics*, 25(2)

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. 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 , and a causal ordering π , the staged tree for this pair 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 $x_1, x_1 x_2, \dots, x_1 \dots x_p$ where $x_i \in \mathcal{X}_{\pi_i}$. Since for the edge $e = ((x_1 \dots x_k), (x_1 \dots x_k x_{k+1}))$ we have $x_{\theta(e)} = \mathbb{P}(x_{k+1} | x_1 \dots x_k)$, 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 their 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 labeled, 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 contain nodes generated from taking 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 $(x_1 \dots x_{k-1}, x_1 \dots x_k)$ represents $P(x_k | x_1 \dots x_{k-1})$, which is also the value of the parameter associated to this edge. Suppose we fix a node $n = a_1 \dots a_k \in L_k$. Each edge coming out of n gives the probabilities for the variable in

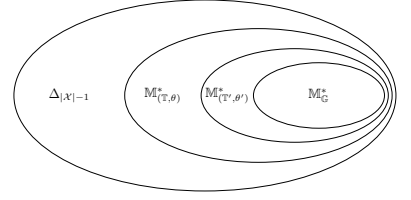


Figure 3.1: Hierarchy of models on \mathcal{X} that we are concerned with, the most general space of models being probability simplex since it contains all distributions on \mathcal{X} , followed by the space of staged tree models $\mathbb{M}_{(\mathbb{T}, \theta)}^*$, the space of CStree models $\mathbb{M}_{(\mathbb{T}', \theta')}^*$ 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 $\mathbb{P}(X_1, \dots, X_p) = \mathbb{P}(X_p | X_{p-1}, \dots, X_1) \mathbb{P}(X_{p-1} | X_{p-2}, \dots, X_1) \dots \mathbb{P}(X_2 | X_1) \mathbb{P}(X_1)$

the next level L_{k+1} , conditioned on the context ($X_{\pi_1} = a_1, \dots, X_{\pi_k} = a_k$). Thus, we can view this node n as representing 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.

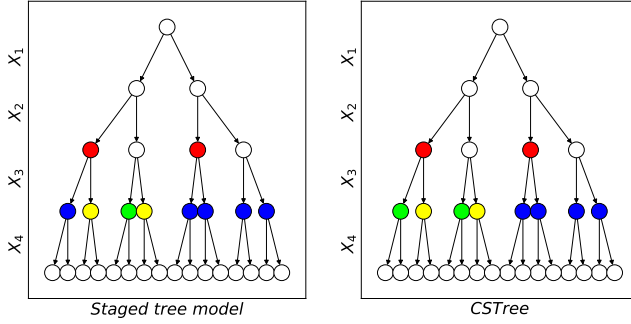


Figure 3.2: Example of a staged tree model that is not a CStree (Left) and a CStree (right) for binary variables X_1, X_2, X_3, X_4 in that causal ordering.

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 1. In this case, the left edge coming out of the root represents $\mathbb{P}(X_1 = 0)$ and the right edge coming out 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 $(X_1 = 1, X_2 = 0, X_3 = 0)$, $(X_1 = 1, X_2 = 1, X_3 = 1)$, $(X_1 = 1, X_2 = 0, X_3 = 1)$, $(X_1 = 1, X_2 = 1, X_3 = 0)$. 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

⁴ Duarte, E. and Solus, L. (2021). Representation of context-specific causal models with observational and interventional data. *arXiv preprint arXiv:2101.09271*

fixed by the stage S_i .

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.

One question that arises is how many CStrees are there for a given system of variables X_1, \dots, X_p taking values in $\mathcal{X}_1, \dots, \mathcal{X}_p$. The authors of the original CStree paper ⁵ provide a closed form expression for the number of CStrees on p binary variables, alongside a table comparing the number of DAGs (sequence A003024 in oeis.org), CStrees and compatibly labeled staged trees, which we include below in Table 3.1. A closed form expression for general CStrees remains to be discovered.

p	DAGs	CStrees	Compatibly labeled staged trees
1	1	1	1
2	3	4	4
3	25	96	180
4	543	59136	2980800
5	29281	26466908160	15619603855888000

⁵ Duarte, E. and Solus, L. (2021). Representation of context-specific causal models with observational and interventional data. *arXiv preprint arXiv:2101.09271*

Table 3.1: Number of DAGs, CStrees and compatibly labeled staged trees on p binary 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 $\pi_1 \dots \pi_p$ 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 ⁶.

Proposition 3.6 (CStrees corresponding to DAGs) *A compatibly labeled staged tree \mathbb{T} with causal ordering $\pi_1 \dots \pi_p$, 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} . ⁷

⁶

⁷ 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 ⁸.

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 i.e. it takes a DAG and returns the corresponding CStree 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 contains 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 contexts in Algorithm 1 which is all the elements of 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

⁸ Tarjan, R. E. (1976). Edge-disjoint spanning trees and depth-first search. *Acta Informatica*, 6(2):171–185

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

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 $X_{\{1,2,3,4\}} = 0110, X_{\{1,2,3,4\}} = 0011$ ⁹ 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.

⁹ $X_{\{1,2,3,4\}} = 0110$ is shorthand for $(X_1 = 0, X_2 = 1, X_3 = 1, X_4 = 0)$

We now describe the algorithm for learning a CStree.¹⁰

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  $\text{COLOUR}(n_1) = \text{COLOUR}(n_2)$  then
8       skip
9     else
10       $c \leftarrow \text{COMMONCONTEXT}(n_1, n_2)$ ;
11       $\text{same\_distr} = \text{TEST}(c, n_1, n_2, D, O[l+1])$ ;
12      if  $\text{same\_distr}$  then
13         $\text{new\_nodes} \leftarrow \text{NODESWITHCONTEXT}(ns, c)$ ;
14         $S \leftarrow \text{UPDATESTAGES}(S, c, \text{new\_nodes})$ ;
15      end
16    end
17  end
18   $l = l + 1$ ;
19 end
20 return  $T, S$ 

```

¹⁰ 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 by assigning them all the same colour.

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 i.e. it merges all possible stages whilst maintaining the CStree property 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 at most 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 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. Since we do learn a DAG model to get possible causal orderings, we can encode the CI relations within this initial DAG into the CStree using Algorithm 1 before learning context-specific information using. However, this could result in CI relations being encoded even though they hold in specific contexts, for example, we might encode $(X_i \perp\!\!\!\perp X_j | X_k)$ even though in reality we only have the CSI relation $(X_i \perp\!\!\!\perp X_j | X_k = 0)$. Thus in the case where we do not know the true causal ordering, we split the algorithm into 2 cases, one where we use the DAG CI relations and another where do not use them, so we can refer them explicitly later on in the experiments section.

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.

Algorithms 3 and 4 consider 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 apriori¹¹, for example a temporal relation between nodes known through physical laws. In this case we can either start statistical testing from an empty tree i.e. use Algorithm 2, or apply the PC algorithm to the data and find a DAG in the Markov Equivalence class with the known ordering so we can encode its CI relations, and then run additional testing to determine context-specific CI relations.

Unlike DAGs, as the number of variables increases 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 2.10 to define a characterization of Markov Equivalence for CStrees¹². We start by first describing the procedure¹³ to gener-

¹¹ Thwaites, P., Smith, J. Q., and Riccomagno, E. (2010). Causal analysis with chain event graphs. *Artificial Intelligence*, 174(12):889–909; and Silander, T. and Leong, T.-Y. (2013). A dynamic programming algorithm for learning chain event graphs. In *International Conference on Discovery Science*, pages 201–216. Springer

¹² Duarte, E. and Solus, L. (2021). Representation of context-specific causal models with observational and interventional data. *arXiv preprint arXiv:2101.09271*

¹³ 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

Algorithm 3: LEARNCSTREEWITHOUTDAGCI

Learning a CStree from observational data when the true causal ordering is unknown and we do not encode the CI relations in the initial DAG

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    $G \leftarrow g \in CPDAG$  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      $orderings \leftarrow ALLTOPOLOGICALSORT(G)$ ;
13   end
14   for  $O$  in  $orderings$  do
15      $T \leftarrow \text{Empty CStree with ordering } O$ ;
16      $T, S \leftarrow LEARNCSTREE(T, S, O, D)$ ;
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  $min\_stage\_trees$ 

```

Algorithm 4: LEARN CSTREE WITH DAGCI

Learning a CStree from observational data when the true causal ordering is unknown and we encode the CI relations in the initial DAG

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      $T, S \leftarrow \text{LEARN CSTREE}(T, S, O, D)$ ;
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  $min\_stage\_trees$ 

```

ate the CSI relations from a CStree and its stages, which uses Lemma 3.5

Algorithm 5: 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 5 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 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 5 can be placed as 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 of set of all CSI relations satisfied by each distribution $\mathbb{P} \in \mathbb{M}_{(\mathbb{T}, \theta)}$ includes the CSI relations recovered from Algorithm 5, 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 $(X_A \perp\!\!\!\perp X_B | X_C = x_C) \in \mathbb{J}$ then $(X_B \perp\!\!\!\perp X_A | X_C = x_C) \in \mathbb{J}$
2. Decomposition, If $(X_A \perp\!\!\!\perp X_{B \cup D} | X_S, X_C = x_C) \in \mathbb{J}$ then $(X_A \perp\!\!\!\perp X_B | X_S, X_C = x_C) \in \mathbb{J}$
3. Weak union, If $(X_A \perp\!\!\!\perp X_{B \cup D} | X_S, X_C = x_C) \in \mathbb{J}$ then $(X_A \perp\!\!\!\perp X_B | X_{S \cup D}, X_C = x_C) \in \mathbb{J}$
4. Contraction, If $(X_A \perp\!\!\!\perp X_B | X_{S \cup D}, X_C = x_C) \in \mathbb{J}$ and $(X_A \perp\!\!\!\perp X_D | X_S, X_C = x_C) \in \mathbb{J}$ then $(X_A \perp\!\!\!\perp X_{B \cup D} | X_S, X_C = x_C) \in \mathbb{J}$
5. Intersection, If $(X_A \perp\!\!\!\perp X_B | X_{S \cup D}, X_C = x_C) \in \mathbb{J}$ and $(X_A \perp\!\!\!\perp X_B | X_{B \cup D}, X_C = x_C) \in \mathbb{J}$ then $(X_A \perp\!\!\!\perp X_{B \cup S} | X_D, X_C = x_C) \in \mathbb{J}$
6. Specialization, If $(X_A \perp\!\!\!\perp X_B | X_S, X_C = x_C) \in \mathbb{J}$ and $T \subset S, x_T \in \mathcal{X}_T$ then $(X_A \perp\!\!\!\perp X_B | X_{S \setminus T}, X_{T \cup C} = x_{T \cup C}) \in \mathbb{J}$
7. Absorption, If $(X_A \perp\!\!\!\perp X_B | X_S, X_C = x_C) \in \mathbb{J}$ and $\exists T \subset C$ such that $\forall x_T \in \mathcal{X}_T$ we have $(X_A \perp\!\!\!\perp X_B | X_S, X_{C \setminus T} = x_{C \setminus T}, X_T = x_T) \in \mathbb{J}$ then $(X_A \perp\!\!\!\perp X_B | X_{S \cup T}, X_{C \setminus T} = x_{C \setminus T}) \in \mathbb{J}$

Given a set of context-specific CI relations \mathbb{J} , the successive application of the axioms above results in the context-specific closure of the \mathbb{J} , denoted $\bar{\mathbb{J}}$. Given a set of variables \mathbb{V} , a context-specific conditional independence model \mathbb{J}^* is a set of quadruples $\langle A, B | S, X_C = x_C \rangle$ where A, B, S, C have no intersection and are subsets of \mathbb{V} , and S, C could be empty, and $\langle \{\}, B | S, X_C = x_C \rangle, \langle A, \{\} | S, X_C = x_C \rangle$ are in \mathbb{J}^* . If \mathbb{J}^* satisfies the context-specific conditional independence axioms above, we call it a context-specific graphoid, which generalizes the notion of graphoids ¹⁴.

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 ¹⁵.

Definition 3.10 (Minimal contexts) *Given a the Context Specific closure $\bar{\mathbb{J}}$ of the CSI relations entailed by a CStree, we say that $X_C = x_C$ is a minimal context if we have $(X_A \perp\!\!\!\perp X_B | X_S, X_C = x_C) \in \bar{\mathbb{J}}$ and there is no non-empty subset $T \subset C$ such that $(X_A \perp\!\!\!\perp X_B | X_{S \cup T}, X_{C \setminus T} = x_{C \setminus T}) \in \bar{\mathbb{J}}$.*

Intuitively, the minimal contexts are the smallest contexts that get left behind after repeated application of the Absorption axiom. By

¹⁴ Pearl, J. and Paz, A. (1986). Graphoids: Graph-based logic for reasoning about relevance relations or when would x tell you more about y if you already know z? In *ECAI*

¹⁵ Duarte, E. and Solus, L. (2021). Representation of context-specific causal models with observational and interventional data. *arXiv preprint arXiv:2101.09271*

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 context-specific conditional axioms are key to represent the CStree as a sequence of DAGs. For now, we denote $\mathcal{C}(\mathbb{T})$ as the set of all minimal contexts of a CStree \mathbb{T} , for each minimal context $(X_C = x_C)_i \in \mathcal{C}(\mathbb{T})$, there exists a DAG that encodes the CI relations that hold under this context. We denote $\mathcal{G}(\mathbb{T}) := \{\mathcal{G}_{X_C=x_C}\}_{X_C=x_C \in \mathcal{C}(\mathbb{T})}$ 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 $(X_A \perp\!\!\!\perp X_B \mid X_S, X_C)$.

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 for all outcomes x_T .¹⁶ 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 6 describes our approach, which makes use of Algorithm 7 as a subroutine to find the largest possible T as in the definition of the absorption axiom.¹⁷

Theorem 3.13 Given a set of Context Specific closure \bar{J} of CSI relations

¹⁶ If we are to use the CSI relations extracted from Algorithm 5 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 $B, C \subset [p-1]$

¹⁷ 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 . `VARIABLESOF-CONTEXT` 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

Algorithm 6: 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     $T\_size \leftarrow \text{RECURSIVETEST}(T\_sizes, C, cs);$ 
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         $possible\_mc \leftarrow [(var, val) \text{ for } (var, val) \text{ in } rel\_context \text{ if } var \text{ not in } possible\_T];$ 
20         $outcome\_count = 0;$ 
21        for  $outcome$  in  $all\_outcomes$  do
22           $contained \leftarrow$ 
23             $\text{CONTAINED}(outcome, cs, possible\_mc);$ 
24          if  $contained = \text{True}$  then
25             $outcome\_count = outcome\_count + 1;$ 
26          end
27          if  $outcome\_count = |all\_outcomes|$  then
28             $T = possible\_T$ 
29          end
30        end
31      end
32     $MCs \leftarrow \text{UPDATEMINIMALCONTEXTS}(MCs, rel, T, J);$ 
33  end
34 end

```

Algorithm 7: 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

```

entailed by a CStree, Algorithm 6 is correct, i.e. it returns the set of minimal contexts.

Proof: We work in sets of CSI relations where the variables X_A, X_B and the conditioning set X_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_r , Algorithm 7 aims to find the largest subset T of C_r 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. This means in order to maximize this size we have to search over the sizes $k_1, \dots, |C_r|$. 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 $|C_r|$ since it is a subset of C_r) 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 $C_r \setminus T_i$ 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 $X_{C_r \setminus T_i} = x_{C_r \setminus T_i}$ as a subcontext, which is consistent with the absorption axiom.

A crude time complexity approximation is as follows. Assume we have p variables. 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 $\mathbb{J}_{A,B,S}$ which share the same variables X_A, X_B and conditioning set X_S . In the worst case, the possible variables for T can be all the variables, or none of them, in which case the algorithm to search for the size of T will run $\log p$ iterations. Within this algorithm, the first midpoint is then $p/2$, which gives a total of $\binom{p}{p/2}$ possible subsets of p with this size. For each of set, we generate the Cartesian product, which has at most $d^{p/2}$ elements, and checking if the required element as per the absorption axiom is in $\mathbb{J}_{A,B,S}$ can take constant time if we use a hash table. Generalizing this, at iteration i , we will look over $\binom{p}{p/2^i}$ possible subsets and corresponding $d^{p/2^i}$ outcomes. Doing this at most $\log p$ times gives the time complexity as

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

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. We note that in practice, all of the above first assumes we have the full set of CSI relations from the context specific closure of the CSI relations entailed by the CStree. This involves repeated application of the axioms above until no new new CSI relations are generated, and in practice is feasible for up to 4 or 5 variables depending on the independence structure of the CStree. This presents itself to be the main bottleneck in this process for the time being.

This motivates us look at the pairwise case, where we focus on Context Specific Independence relations 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 . The definition of pairwise minimal contexts is given below.

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 graphoids¹⁸. An important question is whether the context-specific closure of CSI relations from a CStree form a compositional context-specific graphoid.

Pairwise minimal contexts are always minimal contexts however we

¹⁸ Sadeghi, K. and Lauritzen, S. (2014). Markov Properties for Mixed graphs. *Bernoulli*, 20(2):676 – 696

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 the context-specific causal information contained in 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 that hold under each minimal context can be represented as DAGs, which are called Minimal Context DAGs. Thus, the CStree can be represented as a sequence of DAGs, one for each minimal context. This representation is a consequence of the following theorem ¹⁹.

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 theorem allows us to generate the Minimal Context DAGs once we have the context specific closure of the CSI relations entailed in the CStree. Namely, we take the CI relations that hold under each minimal context, and generate the minimal I-MAP ²⁰, which we can recover with the following procedure ²¹.

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 8 is correct i.e. it returns the minimal context DAGs for each minimal context 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 2.

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*

¹⁹ Duarte, E. and Solus, L. (2021). Representation of context-specific causal models with observational and interventional data. *arXiv preprint arXiv:2101.09271*

²⁰ Verma, T. and Pearl, J. (1990). *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

²¹ Solus, L., Wang, Y., and Uhler, C. (2021). Consistency Guarantees for Greedy Permutation-Based Causal Inference Algorithms. *Biometrika*. asaa104

Algorithm 8: GENERATEMINCONTEXT DAGS

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  $MCDAGS$ 

```

minimal context DAGs $\mathcal{G}(\mathbb{T})$.

To see why we need to consider all Markov equivalent DAGs in the case we do not know the true causal ordering, suppose we have the following case with 4 binary variables.

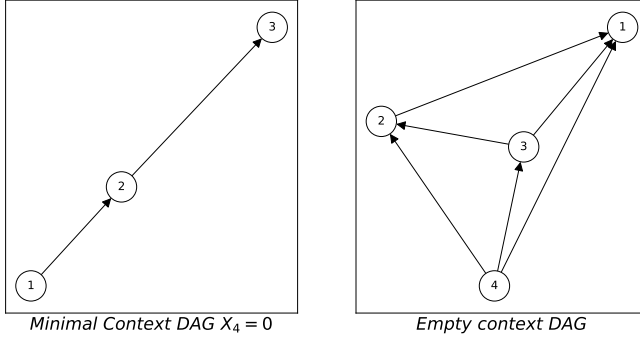


Figure 3.3: Example on why to consider all topological ordering of all Markov equivalent DAGs when learning a CStree

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 exact DAG as our CPDAG from the PC algorithm step, as this DAG only has one topological ordering which is 4321, there is no way to join the stages for the empty context graph with this order to encode $(X_1 \perp\!\!\!\perp X_3 \mid X_4 = 0)$, in which case we do not learn the true model.

We now present the main theorem of this paper, which is the consistency guarantee for learning the CStree when we know the true causal ordering.

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 2 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 $\mathcal{G}_{\mathbb{T}}$. By Theorem 3.15, this is equivalent to \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 $\mathcal{G}_{\mathbb{T}}$. 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

If we do not know the true causal ordering, 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 razor ²² since few stages imply lower model complexity and all else being equal, 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 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 ²³.

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

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 is a multi-dimensional array ²⁴ 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

²² Pearl, J. (2009). *Causality*. Cambridge University Press

²³ Here we compactify the notation further - $X_{\{1, \dots, p\}} = x_{\{1, \dots, p\}}$ is the same as $X_{\{1, \dots, p\}} = x_1 \dots x_p$ which is $X_1 = x_1, \dots, X_p = x_p$.

²⁴ This is the same as tensors in the context of computer science

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 ²⁵ for a tree with levels $(L_1, \dots, L_p) \sim (X_1, \dots, X_p)$.

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

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 ²⁶. The BIC score for a CStree \mathbb{T} with observed data \mathbb{D} is then

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

²⁵ Duarte, E. and Solus, L. (2021). Representation of context-specific causal models with observational and interventional data. *arXiv preprint arXiv:2101.09271*

²⁶ Duarte, E. and Solus, L. (2021). Representation of context-specific causal models with observational and interventional data. *arXiv preprint arXiv:2101.09271*

4

Experiments

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

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 to the method used to determine whether nodes belong to the same stage?
4. How sensitive are the results when using different DAG model learning algorithms to learn the set of possible causal orderings?
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. For minimal contexts, we compute the general minimal contexts rather than the pairwise minimal contexts. We note that for the experiments we run below, the pairwise minimal contexts and general minimal contexts are equivalent.

4.1 Synthetic data

We start by generating random DAGs and generating the corresponding CStrees as per Algorithm 1. We run a sanity check by generating CStrees from DAGs in the 2 extreme cases, the fully connected DAG and the empty DAG. The results are included in the Appendix. For

random DAGs, we first choose the number of variables (nodes in the DAG) p , then we choose an edge with probability $p_{edge} \in (0, 1)$, and keep the 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 8.

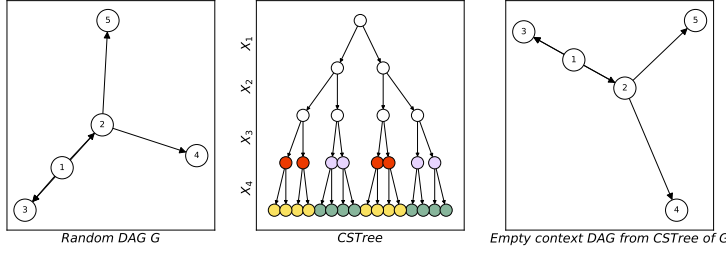


Figure 4.1: Generating CStrees from random DAGs and recovering the original random DAG from the CStree with ordering 1234 using Algorithm 8

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

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

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 4 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 $X_{(PA_{X_i})_1}, \dots, X_{(PA_{X_i})_d}$, 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 simulate samples on 4 binary variables using the above procedure

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. We see the space complexity increase exponentially as mentioned. The time complexity however can depend on the DAG structure, however the trend suggests an exponential increase.

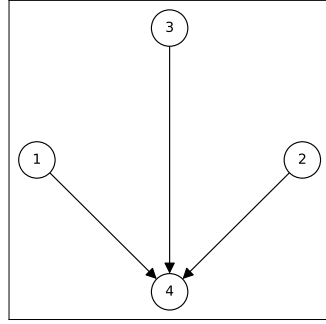


Figure 4.2: Random DAG generated to test Algorithm 3 and 8 on synthetic data. The exact probability distribution values are

$$\begin{aligned}
 &P(X_4 = 0 \mid X_{\{1,2,3\}} = 000) = 0.06, \\
 &P(X_4 = 0 \mid X_{\{1,2,3\}} = 001) = 0.85, \\
 &P(X_4 = 0 \mid X_{\{1,2,3\}} = 010) = 0.87, \\
 &P(X_4 = 0 \mid X_{\{1,2,3\}} = 011) = 0.1, \\
 &P(X_4 = 0 \mid X_{\{1,2,3\}} = 100) = 0.1, \\
 &P(X_4 = 0 \mid X_{\{1,2,3\}} = 101) = 0.16, \\
 &P(X_4 = 0 \mid X_{\{1,2,3\}} = 110) = 0.94, \\
 &P(X_4 = 0 \mid X_{\{1,2,3\}} = 111) = 0.86, \\
 &P(X_1 = 0) = 0.14, P(X_2 = 0) = 0.39, P(X_3 = 0) = 0.37
 \end{aligned}$$

by fixing a DAG and a distribution, and sampling 2000 values from it. The random DAG we generate is shown in Figure 4.2 We then learn the CStree using Algorithm 3 and recover the minimal context DAGs via Algorithm 8 for a random subsample of 500 and all 2000 of these samples. The results of these are shown in Figure We choose 4 variables due to limitations of applying the context specific conditional independence axioms. We also set the ordering to be 1234.

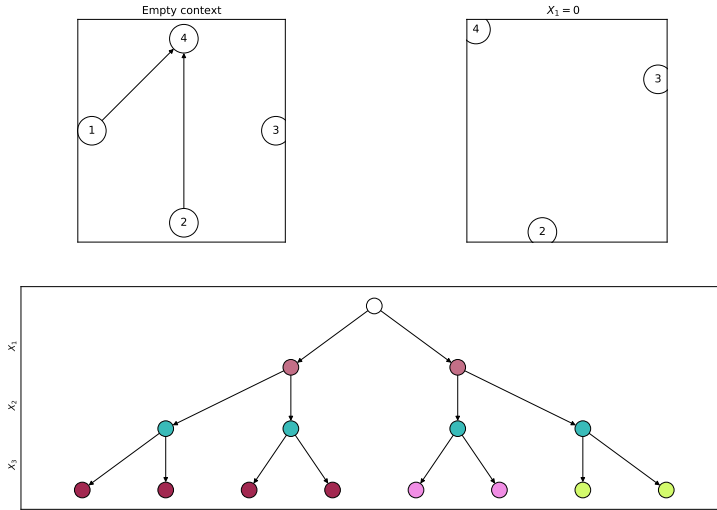


Figure 4.3: CStree and corresponding minimal context DAGs from the 500 samples from the synthetic DAG experiment.

We expect that in the limit of large data, the CSI relations we learn would converge to the CI relations implied by the DAG. 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. Here we use the Epps-Singleton test ¹. We see that the empty context DAG does not

¹ Epps, T. and Singleton, K. J. (1986). An omnibus test for the two-sample problem using the empirical characteristic function. *Journal of Statistical Computation and Simulation*, 26(3-4):177-203

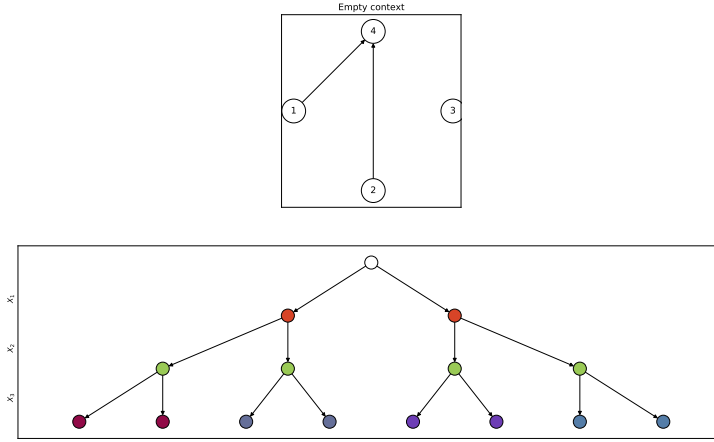


Figure 4.4: CStree and corresponding minimal context DAG from the 2000 samples from the synthetic DAG experiment.

correspond to the true DAG in this case. Although this did not happen significantly often on our synthetic data experiments, in practice the CI testing used to learn the DAG model and the CSI relations can be highly unstable, and this will be a common theme in this chapter. For example it is possible for the CI testing to erroneously encode a CI relation, implying that that the independence relation holds for all outcomes whilst in reality it only holds under specific contexts. Similarly, it is possible for the CSI tests to erroneously merge stages, entailing a CSI relation that does not exist.

In order to assess the influence of the CI relations from the DAG in the final CStree, we tested the use of Algorithm 3 on the same 2000 samples used to get the results in 4.4. The results are shown in Figure 4.5. In this case we do recover the true DAG. Since this CStree was learnt only using the CSI relations, it is reasonable to assume that the errors in CI testing contributed to the previous result. This suggests that we could expect better performance if we only use DAG discovery algorithms, such as the PC algorithm, to learn possible orderings of the true CStree rather than to learn an approximation of the CStree model itself through the CI relations encoded in them.

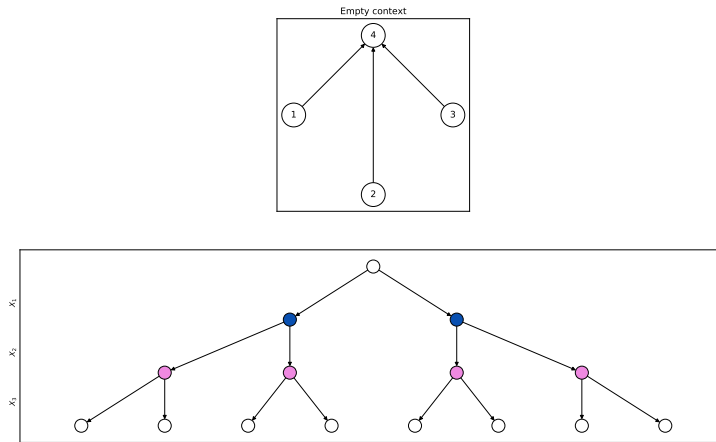


Figure 4.5: CStree and corresponding minimal context DAG from the same samples as in Figure 4.4 but without encoding the CI relations from the initial DAG.

4.2 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 ². 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}

This data set does not require further pre-processing since the outcome 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 ³. 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, from Algorithm 1.
2. The CStree learnt only from context specific information without encoding the CI relations from the DAG, from Algorithm 3.
3. The CStree learnt from encoding the CI relations from the DAG and then learning further context specific information, from Algorithm 4.

² Reinis, Z., Pokorny, J., Basika, V., Tiserova, J., Gorican, K., Horakova, D., Stuchlikova, E., Havranek, T., and Hrabovsky, F. (1981). Prognostic significance of the risk profile in the prevention of coronary heart disease. *Bratisl Lek Listy*, 76:137–150

³ Koller, D. and Friedman, N. (2009). *Probabilistic graphical models: principles and techniques*. MIT press; and Ankan, A. and Panda, A. (2015). pgmpy: Probabilistic graphical models using python. In *Proceedings of the 14th Python in Science Conference (SCIPY 2015)*. Citeseer

For each of the cases above we create/merge stages using the Anderson-Darling test ⁴, the Epps-Singleton test ⁵, and the symmetric KL divergence ⁶ with different threshold values. The use of the symmetric KL divergence is motivated by the unreliability of conditional independence testing in practice ⁷. Given 2 distributions \mathbb{P}, \mathbb{Q} over the same discrete state space \mathcal{X} , the symmetric KL divergence is as follows.

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

The symmetric KL (SKL) divergence is a measure of distance between the distributions \mathbb{P} and \mathbb{Q} , and returns 0 if $\mathbb{P} = \mathbb{Q}$ for all outcomes $x \in \mathcal{X}$.

The output of the PC algorithm is a CPDAG, and whilst the outcome of the Hill climbing algorithm is a DAG, there are computationally efficient algorithms that allow us to generate the CPDAG that represents its MEC ⁸. Once we have the CPDAG, generating the DAGs consistent with it involves orienting the undirected edges such that we do not form any new v-structures and do not create cycles.

We apply Algorithms 1, 3, and 4 to the setup above, and record the CStrees with the minimum number of stages over all possible causal orderings, alongside their BIC scores. Note that it is possible for multiple non-statistically (non-Markov) equivalent CStrees to have stages equal to the minimum number of stages. This can occur, for example when the underlying true causal structure is not faithful to a CStree but is instead faithful to a non-CStree staged tree. It can also occur due to testing errors induced by small sample size. In this case, we select the BIC-optimal CStree since the BIC is a score-equivalent and consistent estimator for CStrees ⁹. The score-equivalence property implies in the limit of large data, all statistically equivalent equivalent CStrees will have the same BIC score. We call this **experimental setup 1**, and the results of this experiment for the coronary dataset are summarized in the Table 4.2.

The CPDAG from the PC algorithm gave 12 orderings, whilst the Hill climbing algorithm provided 16, where we optimized for the BIC score.

Here we see that the minimum number of stages over all the configurations is 6, which results in the CStree shown in Figure 4.6. The causal ordering is 346125, and we get this CStree by learning the original DAG using the PC algorithm, followed by the Anderson-Darling

⁴ Scholz, F. W. and Stephens, M. A. (1987). K-sample anderson-darling tests. *Journal of the American Statistical Association*, 82(399):918–924

⁵ Epps, T. and Singleton, K. J. (1986). An omnibus test for the two-sample problem using the empirical characteristic function. *Journal of Statistical Computation and Simulation*, 26(3-4):177–203

⁶ Also called the Kullback-Leibler divergence

⁷ Shah, R. D. and Peters, J. (2020). The hardness of conditional independence testing and the generalised covariance measure. *The Annals of Statistics*, 48(3)

⁸ Chickering, D. M. (2002a). Learning equivalence classes of bayesian-network structures. *J. Mach. Learn. Res.*, 2:445–498

⁹ Duarte, E. and Solus, L. (2021). Representation of context-specific causal models with observational and interventional data. *arXiv preprint arXiv:2101.09271*

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/o DAG	6	-6793.7	13	-6812.1	23	-6757.1	31	-7032.3	44	-6753.4
	CSTree with DAG	8	-6771.3	11	-6780.4	15	-6775.4	18	-6739.4	18	-6739.4
HC	DAG	28	-6714.78	28	-6714.8	28	-6714.8	28	-6714.8	28	-6714.8
	CSTree w/o DAG	7	-6791.0	30	-6841.9	24	-6802.9	36	-6828.1	47	-6755.9
	CSTree with DAG	8	-6788.7	15	-6809.7	21	-6794.9	20	-6786.9	27	-6711.0

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 CI relation(s) encoded in the DAG may have put nodes in the same stage, which otherwise may have been merged from the context specific testing procedure. These observations collectively imply that, possibly due to the small sample size, or violation of the CStree faithfulness assumption, the minimum stage condition for selecting the best CStree is not reasonable. Instead, we can take the BIC-optimal CStree with the fewest stages learned over all the configurations, as this provides a balance between the model fit and sparsity.

We note that we do learn a CStree with 8 stages which is not a DAG model, corresponding to the ordering 431265, which we show below in Figure 4.7.

Among the trees with the fewest stages, whose information is shown in Table 4.2, the CStree with the highest BIC score we recorded is generated with the ordering 541326, using the symmetric KL

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.

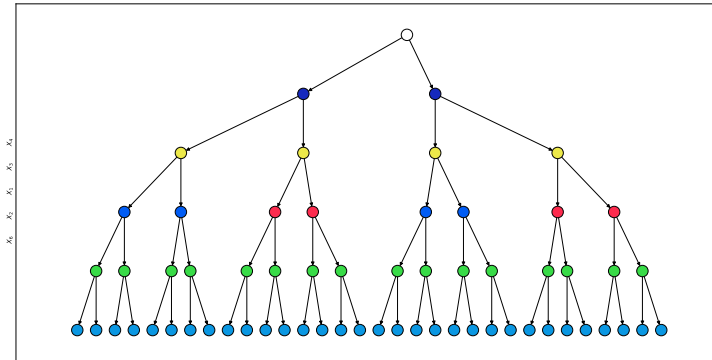


Figure 4.6: CStree for the coronary dataset with the lowest number of stages, which has ordering 346125 and corresponds to a DAG model.

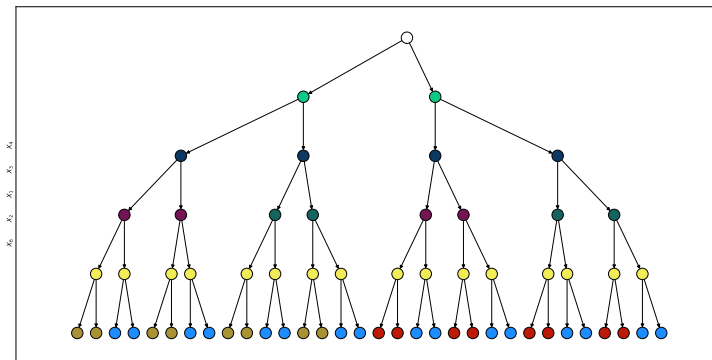


Figure 4.7: CStree with causal ordering 431265 learnt from the coronary dataset which has 8 stages and does not correspond to a DAG model.

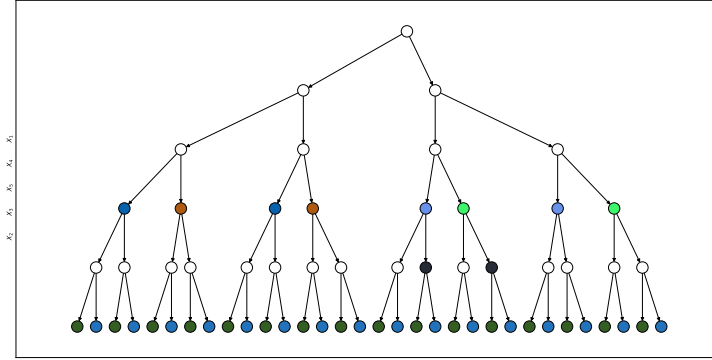


Figure 4.8: CStree with the highest BIC score among all CStrees with the minimum number of stages in experimental setup 1 for the coronary dataset, which has ordering 541326.

divergence with a threshold of 5×10^{-7} 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 30 stages, which motivates us to define the **experimental setup 2**, where we find the CStree with the 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.

The results of Table 4.2 and 4.3 show that the BIC-optimal CStree configuration still remains the same. In general, the differences between experimental setup 1 and 2 on this dataset show that we can find CStrees with a higher BIC score given the same configuration i.e. merging procedure and initial DAG learning algorithm, at the cost of more stages.

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 the BIC score ¹⁰. 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.
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

¹⁰ Carli, F., Leonelli, M., Riccomagno, E., and Varando, G. (2020). The r package stagedtrees for structural learning of stratified staged trees. *arXiv preprint arXiv:2004.06459*

Goal: Max BIC		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/o DAG	7	-6783.7	13	-6812.1	23	-6757.1	43	-6749.6	44	-6753.4
	CSTree with DAG	8	-6771.3	11	-6780.4	15	-6775.4	18	-6739.4	18	-6739.4
HC	DAG	31	-6714.8	28	-6714.8	28	-6714.8	28	-6714.8	28	-6714.8
	CSTree w/o DAG	8	-6775.5	31	-6826.4	27	-6802.8	39	-6828.0	47	-6755.9
	CSTree with DAG	9	-6773.2	16	-6794.1	24	-6786.6	20	-6786.49	27	-6711.0

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 3/4 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 in Table 4.4. None of the staged trees learnt by algorithms (1)-(4) above are CStrees, and they are included in the Appendix.

We see that the staged trees achieve a higher BIC score, which is probably due to their greater flexibility.

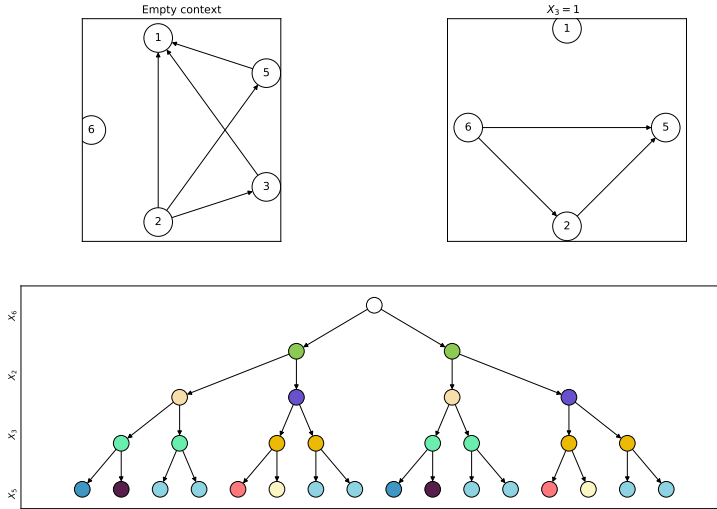
In order to generate the sequence of DAGs that capture the context specific causal information from this dataset, we remove one variable for computational feasibility. We remove the variable that results

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.

Model	Algorithm	BIC	Stages
CStree	Modified Algorithm 3 / 4	-6729.5	30
Non-CStree	Hill climbing	-6645.3	11
Non-Cstree	Backward hill climbing	-6641.1	13
Non-CStree	Backward joining	-6790.4	57
Non-CStree	Hierarchial clustering	-6654.7	10

Table 4.4: Experimental setup 3 on the coronary dataset

in the highest BIC score after its removal, which in this case is the variable corresponding to systolic blood pressure (variable 4). The initial DAG is generated from the PC algorithm, and we learn the tree structure using both the CI relations from the DAG and the context specific staging procedure using the symmetric KL divergence with a threshold of 5×10^{-5} . We choose the BIC optimal tree rather the one with the fewest stages, since the latter results in a context specific closure which could be computed in a reasonable amount of time. This results in the CStree with ordering 62351 which has 10 stages, shown in Figure 4.9.

Figure 4.9: BIC-optimal CStree alongside the minimal context DAGs on the coronary dataset, which has the ordering 62351. The initial DAG is learnt with the PC algorithm, and the stages were merged using the symmetric KL divergence with a threshold of 5×10^{-5} .

We note that in some cases one can read off the minimal contexts straight from the CStrees. One example is the CStree in Figure 4.8, where we can see that the minimal contexts are the empty context and $(X_5 = 1, X_4 = 0, X_3 = 1)$, where one can draw the minimal I-MAPs for the whole tree and the one for just this context through inspecting the CStree. Thus in some cases, despite the computational limitations of finding the minimal contexts automatically, we still see that we can use CStrees to more easily read context-specific causal information. In particular, it is much easier to interpret the minimal

context DAGs in Figure 4.9 than its corresponding CStree which we learnt from observed data.

4.3 Mice protein expression data

This is a dataset with expression levels of 77 proteins, measured in the cerebral cortex of 8 classes of mice ¹¹. 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 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 ¹² 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 response 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 if we use the PC algorithm for the initial DAG. The chosen variables 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 .

We show results of applying experimental setup 1 on this dataset in Table 4.5 below. The PC algorithm resulted in 196 possible causal

¹¹ Higuera, C., Gardiner, K. J., and Cios, K. J. (2015). Self-organizing feature maps identify proteins critical to learning in a mouse model of down syndrome. *PLOS ONE*, 10(6):e0129126

¹² Guyon, I., Weston, J., Barnhill, S., and Vapnik, V. (2002). *Machine Learning*, 46(1/3):389–422

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/o DAG	7	-6753.0	8	-6572.7	61	-5917.8	64	-5823.9	78	-5807.6
	CStree with DAG	11	-6591.4	11	-6591.4	67	-5702.0	68	-5705.4	71	-5692.2
HC	DAG	502	-6472.1	502	-6472.1	502	-6472.1	502	-6472.1	502	-6472.1
	CStree w/o DAG	8	-6191.5	17	-6068.9	67	-6386.9	90	-6348.8	90	-6348.8
	CStree with DAG	9	-6173.3	9	-6173.3	29	-6015.7	37	-6429.9	45	-6323.8

Table 4.5: Experimental setup 1 on the mice protein expression dataset

orderings, whilst the Hill climbing algorithm resulted in 12. We use the K2 score¹³ for the Hill climbing since otherwise the number of possible orderings increases to 6408.

¹³ Carvalho, A. M. (2009). Scoring functions for learning bayesian networks

The results from Table 4.5 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. The minimum number of stages we get is 7, which is not helpful considering that this is the full independence model. We also note the high number of stages in the DAG model from the Hill climbing algorithm, which is due to the high degree of connectivity in the graph, alongside the variable corresponding to the response class taking 8 possible outcomes. The CStrees corresponding to 8 and 9 stages are also DAG models, however we do learn CStrees with 10 stages which are not DAG models. We show one of these trees below in Figure 4.10.

We show the result of experimental setup 2 on this dataset in Table 4.6.

The CStree with the highest BIC score has 133 stages, with a score of -5433.0. This tree has an ordering of 32765148, and we show it in Figure 4.11

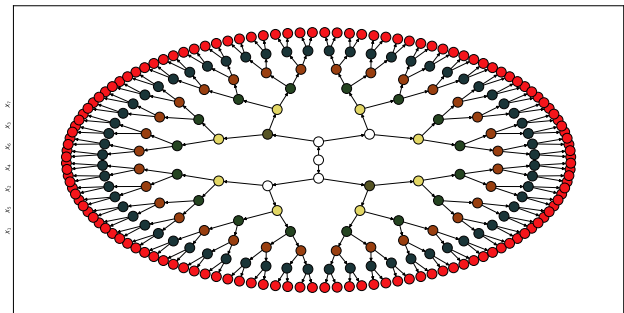


Figure 4.10: CStree with ordering 73642518 learnt on the mice protein expression dataset, and has 10 stages. This CStree does not correspond to a DAG model. The initial DAG was learnt using the PC algorithm, and we use the Epps-Singleton test to merge stages.

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 with DAG	11	-6591.4	36	-5906.2	75	-5598.9	76	-5602.4	79	-5589.1
HC	DAG	502	-6472.1	502	-6472.1	502	-6472.1	502	-6472.1	502	-6472.1
	CSTree w DAG	11	-6067.4	20	-5944.9	68	-5906.6	97	-5679.8	96	-5684.9
	CSTree with DAG	12	-6049.3	25	-5944.8	45	-5993.4	47	-5843.9	55	-5786.0

Table 4.6: Experimental setup 2 on the mice protein expression dataset.

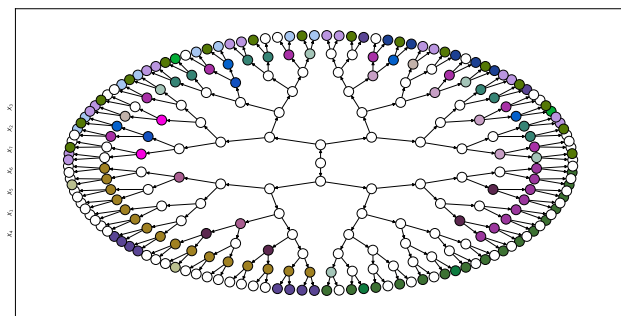


Figure 4.11: CStree with ordering 32765148 learnt on the mice protein expression dataset, and has 133 stages. This CStree does not correspond to a DAG model, and is the BIC optimal tree we learnt on this dataset. The initial DAG was learnt using the PC algorithm, and we use the symmetric KL divergence with a threshold of 5×10^{-6} to merge stages.

We now take the recorded ordering which gave the highest BIC score on this dataset (which is 3276514) and perform experimental setup 3. The results are summarized in Table 4.7.

Model	Algorithm	BIC	Stages
CStree	Modified Algorithm 3 / 4	-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	Hierarchical clustering	-5990.4	17

Table 4.7: Experimental setup 3 on the mice protein expression dataset

We see that for this dataset our BIC optimal CStree has a higher number of stages, whilst more sparse staged trees have a higher BIC score. This could be due to the higher flexibility in how staging could be done in staged trees.

In order to visualize the minimal context DAGs for the mice protein expression dataset, we first again perform recursive feature elimination to get 4 features plus the response variable. This results in the selection of variables 2,3,4,5,8 from the original set of variables. We then apply Algorithm 4 to learn the CStree using the PC algorithm alongside the symmetric KL divergence with a threshold of 5×10^{-6} . The results are shown below in Figure 4.12 below.

We note that all of the CStreets we believe might be good candidate models and thus shown above have a causal ordering ending with 8, which is the response variable for this dataset. This aligns with the aim of this dataset is to discriminate between the 8 classes for this variable, since identifying features that discriminate the response class has the same meaning as finding the features that cause changes in the response class, in which case we would expect the feature values to happen before the response variable values.

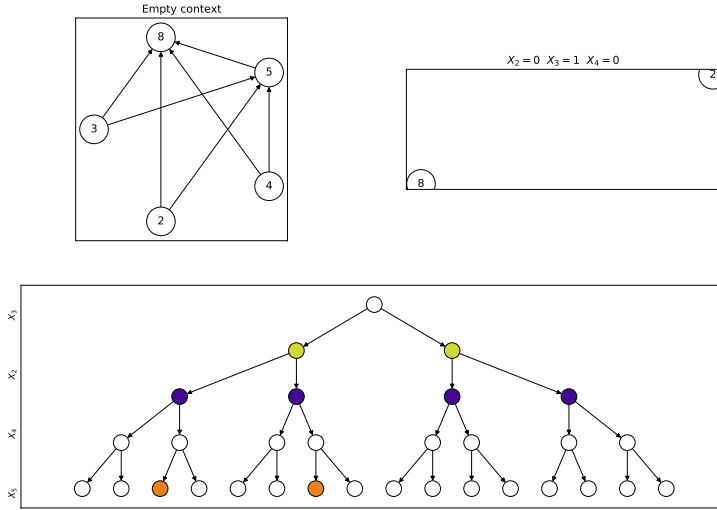


Figure 4.12: Minimum stage CStree learnt from the mice protein expression dataset on the subset of variables 2,3,4,5,8. This CStree has the ordering 32458. The original DAG was learnt using the PC algorithm, and we use the symmetric KL divergence with a threshold of 5×10^{-6} to merge stages.

4.4 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 features¹⁴. 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. Similar to before, we categorize each variable in binary classes depending on whether or not they have higher or lower than the median value of the corresponding samples.

We see that the resulting CStree perhaps surprisingly does not include any non-singleton stages.

¹⁴ Baldi, P., Sadowski, P., and Whiteson, D. (2014). Searching for exotic particles in high-energy physics with deep learning. *Nature Communications*, 5(1)

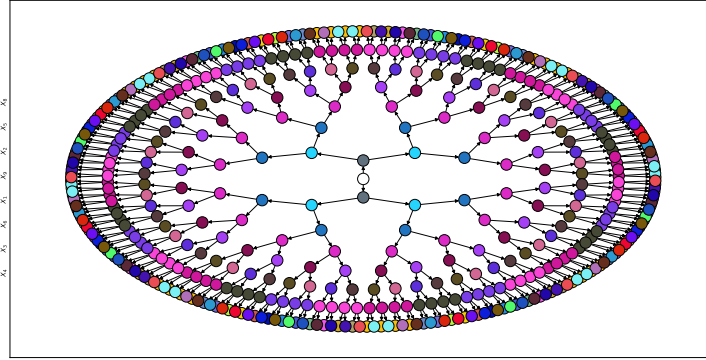


Figure 4.13: First CStree learnt from the supersymmetry dataset, with causal ordering 852916347.

4.5 Vitamin-D data

This dataset comes from a real study on Vitamin D and mortality¹⁵, 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 response variable which is a binary value indicating knowledge of whether the subject passed away during the follow up. There are 2571 samples, with no missing data. 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 .

The results of experimental setup 1 and 2 are shown in Tables 4.8 and 4.9 respectively. We note that the PC algorithm gave us 40 possible causal orderings, whilst the Hill climbing algorithm gave 10.

We again get the full independence model when using the Anderson test. The symmetric KL divergence threshold of 5×10^{-1} also results in the full independence model which suggests that it is too high.

¹⁵ Martinussen, T., Sørensen, D. N., and Vansteelandt, S. (2017). Instrumental variables estimation under a structural cox model. *Biostatistics*, 20(1):65–79

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	-8805.6
	CSTree with DAG	6	-8480.1	12	-8410.1	4	-9041.5	9	-8405.2	12	-8410.1
HC	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	-8805.6
	CSTree with DAG	6	-8480.1	12	-8410.1	4	-9041.5	9	-8405.2	12	-8410.1

Table 4.8: Experimental setup 1 on the Vitamin D dataset

Among the CStrees with the minimum stages, the highest BIC score is -8405.2 which is achieved by using the CI relations learnt from the PC algorithm alongside context specific merging using the symmetric KL divergence with a threshold of 5×10^{-2} . We show the resulting CStree and minimal context DAGs in Figure 4.14 below.

We show the results of experimental setup 2 on the Vitamin D dataset below in Table 4.9. The results are exactly the same, namely, the CPDAG learnt from both the PC algorithm and the Hill climbing algorithm coincide, alongside the entries corresponding to the minimum number of stages and maximum BIC over all orderings per configuration.

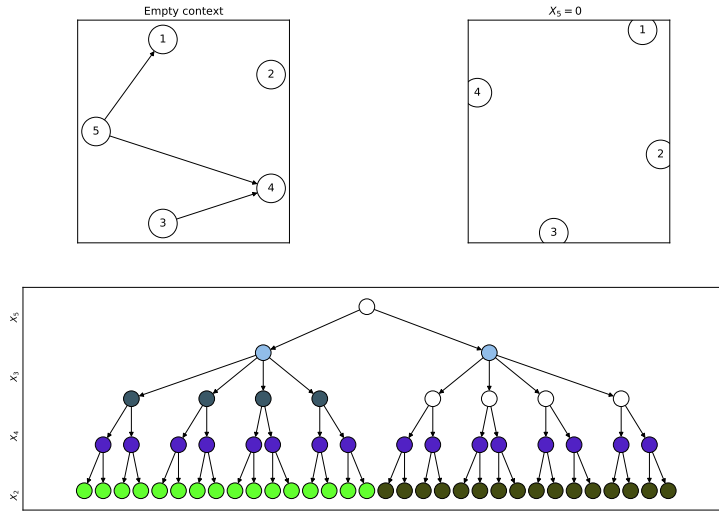


Figure 4.14: CStree and corresponding minimal context DAGs with the highest BIC in both experimental setup 1 and 2 on the Vitamin D dataset. The causal ordering is 53421.

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 with DAG	6	-8480.1	12	-8410.1	5	-8749.9	9	-8405.2	12	-8410.1
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 with DAG	6	-8480.1	12	-8410.1	5	-8749.9	9	-8405.2	12	-8410.1

Table 4.9: Experimental setup 2 on the Vitamin D dataset

Table 4.10 shows the results of experimental setup 3 on the Vitamin D dataset.

Model	Algorithm	BIC	Stages
CStree	Modified Algorithm 3/4	-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	Hierarchical clustering	-8910.0	10

Table 4.10: Experimental setup 3 on Vitamin D dataset with ordering 53421

The results of Table 4.10 show that the CStree does better than all staged tree models learnt in terms of the BIC score, and also provides the sparsest model learnt directly from the data, which is a fair comparison if we exclude the number of stages learnt from 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 result of this is that 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. A balance is achieved when taking the causal ordering 32145, which has a slightly worse BIC score of -8776.5 and 9 stages. As a result, we use this as the basis for the comparison with staged trees. The results are summarized in Table 4.11 below.

We see that in terms of the BIC score the CStree is not optimal, how-

Model	Algorithm	BIC	Stages
CStree	Modified Algorithm 3 / 4	-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 4.11: Experimental setup 3 on Vitamind D dataset with orrdering 32145

ever the model has fewer stages when compared to the staged trees. Despite this, the CStree maintains the benefit of the more intuitive representation in terms of minimal context DAGs. We show these alongside the corresponding CStree in Figure 4.15 below.

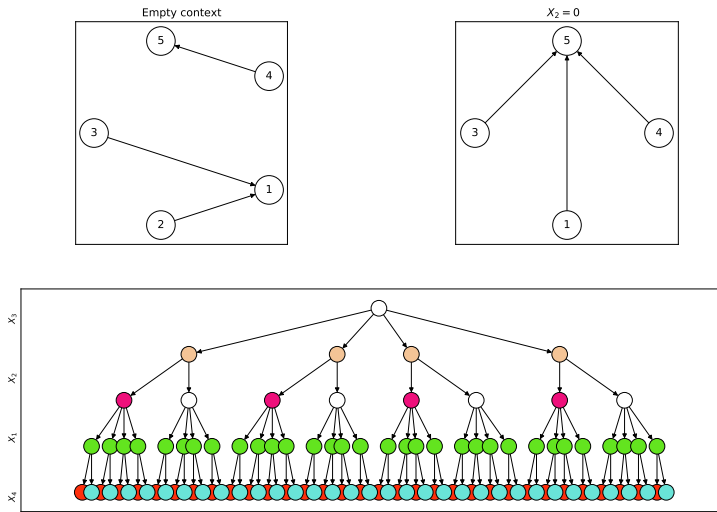


Figure 4.15: CStree and minimal context DAGs of the Vitamin D dataset for the causal ordering 32145.

We see that the minimal context DAGs contain v-structures, meaning the directed edges are fixed in the Markov Equivalence class. Thus under our assumptions, they represent direct causal relationships learnt from the data.

4.6 Discussion of empirical performance

We see in general that using the symmetric KL divergence to identify stages in the CStree result in 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 ¹⁶. We also see that the Anderson-Darling test tends to merge stages more frequently when

¹⁶ Shah, R. D. and Peters, J. (2020). The hardness of conditional independence testing and the generalised covariance measure. *The Annals of Statistics*, 48(3)

compared to other staging procedures, which sometimes result in the full independence model.

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 response 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 orderings, which is a model selection problem. In Algorithms 3 and 4 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 razor¹⁷ - 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.

¹⁷ Pearl, J. (2009). *Causality*. Cambridge University Press

All of the above analysis also takes for granted the validity of the approaches taken to group data into categories in the first place. The approach we take here is on the more simpler side, since we partition variables whos values are ordered. Even if we are given categorical data, it might be in the best interest to group since it is possible to have 2 or more features which provide the same level of information. This would also help with potential issues small data sizes, since more outcomes can thin out samples which have fixed contexts.

If we truly want to optimize for a score in practice, it is inevitable that there will be some level of hyperparameter tuning involved. Some examples include the threshold for the symmetric KL di-

vergence and the value corresponding to the minimum number of increase required to not stop the Hill climbing algorithm. One method is to use a grid search, however this is not scalable, in which case one can use probabilistic search methods to search for a good hyperparameter setting ¹⁸.

¹⁸ Oh, C., Tomczak, J. M., Gavves, E., and Welling, M. (2019). Combinatorial bayesian optimization using the graph cartesian product. *arXiv preprint arXiv:1902.00448*

Conclusions

5.1 Summary

We started 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 covered the limitations of CI relations in comparison to CSI relations, and covered CStrees as a means to encode such CSI relations. We then showed how one can learn these CSI relations from observational data to 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. Based on the empirical results presented in this work, it might also be interesting to jointly optimize for the minimization of the stages and the maximization of some score like the BIC score. 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

¹ Yang, K., Katcoff, A., and Uhler, C. (2018). Characterizing and learning equivalence classes of causal DAGs under interventions. In Dy, J. and Krause, A., editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5541–5550. PMLR

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

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

⁴ Tu, R., Zhang, C., Ackermann, P., Mohan, K., Kjellström, H., and Zhang, K. (2019). Causal discovery in the presence of missing data. In Chaudhuri, K. and Sugiyama, M., 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

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 the combinatorial problem into a purely continuous problem ⁵.

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

Bibliography

- Andersson, S. A., Madigan, D., and Perlman, M. D. (1997). A characterization of markov equivalence classes for acyclic digraphs. *The Annals of Statistics*, 25(2).
- Ankan, A. and Panda, A. (2015). pgmpy: Probabilistic graphical models using python. In *Proceedings of the 14th Python in Science Conference (SCIPY 2015)*. Citeseer.
- Baldi, P., Sadowski, P., and Whiteson, D. (2014). Searching for exotic particles in high-energy physics with deep learning. *Nature Communications*, 5(1).
- Carli, F., Leonelli, M., Riccomagno, E., and Varando, G. (2020). The r package stagedtrees for structural learning of stratified staged trees. *arXiv preprint arXiv:2004.06459*.
- Carvalho, A. M. (2009). Scoring functions for learning bayesian networks.
- Castelletti, F. (2020). Bayesian model selection of gaussian directed acyclic graph structures. *International Statistical Review*, 88(3):752–775.
- Chickering, D. M. (2002a). Learning equivalence classes of bayesian-network structures. *J. Mach. Learn. Res.*, 2:445–498.
- Chickering, D. M. (2002b). Optimal structure identification with greedy search. *Journal of machine learning research*, 3(Nov):507–554.
- Collazo, R. A., Görgen, C., and Smith, J. Q. (2018). *Chain event graphs*. CRC Press.
- Duarte, E. and Solus, L. (2020). Algebraic geometry of discrete interventional models. *arXiv preprint arXiv:2012.03593*.

- Duarte, E. and Solus, L. (2021). Representation of context-specific causal models with observational and interventional data. *arXiv preprint arXiv:2101.09271*.
- Elias Bareinboim, Juan D. Correa, D. I. T. I. (2020). On pearl's hierarchy and the foundations of causal inference.
- Epps, T. and Singleton, K. J. (1986). An omnibus test for the two-sample problem using the empirical characteristic function. *Journal of Statistical Computation and Simulation*, 26(3-4):177–203.
- Fisher, R. (1935). *The design of experiments*. 1935. Oliver and Boyd, Edinburgh.
- Geiger, D., Verma, T., and Pearl, J. (1990). Identifying independence in bayesian networks. *Networks*, 20(5):507–534.
- Glymour, C., Zhang, K., and Spirtes, P. (2019). Review of causal discovery methods based on graphical models. *Frontiers in Genetics*, 10(nil):nil.
- Guyon, I., Weston, J., Barnhill, S., and Vapnik, V. (2002). *Machine Learning*, 46(1/3):389–422.
- Hagberg, A. A., Schult, D. A., and Swart, P. J. (2008). Exploring network structure, dynamics, and function using networkx. In Varoquaux, G., Vaught, T., and Millman, J., editors, *Proceedings of the 7th Python in Science Conference*, pages 11 – 15, Pasadena, CA USA.
- Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M. H., Brett, M., Haldane, A., del Río, J. F., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K., Reddy, T., Weckesser, W., Abbasi, H., Gohlke, C., and Oliphant, T. E. (2020). Array programming with NumPy. *Nature*, 585(7825):357–362.
- Higuera, C., Gardiner, K. J., and Cios, K. J. (2015). Self-organizing feature maps identify proteins critical to learning in a mouse model of down syndrome. *PLOS ONE*, 10(6):e0129126.
- Hu, P., Jiao, R., Jin, L., and Xiong, M. (2018). Application of causal inference to genomic analysis: Advances in methodology. *Frontiers in Genetics*, 9:238.
- Hume, D. (1748). *Enquiry Concerning Human Understanding*.
- Hünermund, P. and Bareinboim, E. (2019). Causal inference and data-fusion in econometrics. *arXiv preprint arXiv:1912.09104*.

- Kalisch, M. and Bühlmann, P. (2007). Estimating high-dimensional directed acyclic graphs with the pc-algorithm. *Journal of Machine Learning Research*, 8(22):613–636.
- Koller, D. and Friedman, N. (2009). *Probabilistic graphical models: principles and techniques*. MIT press.
- Krekel, H., Oliveira, B., Pfannschmidt, R., Bruynooghe, F., Laughner, B., and Bruhin, F. (2004). *pytest* x.y.
- Lauritzen, S. L. (1996). *Graphical models*, volume 17. Clarendon Press.
- MacKay, D. J. C. (1992). *Bayesian Interpolation*, pages 39–66. Springer Netherlands, Dordrecht.
- Martinussen, T., Sørensen, D. N., and Vansteelandt, S. (2017). Instrumental variables estimation under a structural cox model. *Biostatistics*, 20(1):65–79.
- Meek, C. (1995a). Causal inference and causal explanation with background knowledge. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence, UAI’95*, pages 403–410, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Meek, C. (1995b). Strong completeness and faithfulness in bayesian networks. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence, UAI’95*, pages 411–418, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Oh, C., Tomczak, J. M., Gavves, E., and Welling, M. (2019). Combinatorial bayesian optimization using the graph cartesian product. *arXiv preprint arXiv:1902.00448*.
- Pearl, J. (2009). *Causality*. Cambridge University Press.
- Pearl, J. and Mackenzie, D. (2018). *The Book of Why: The New Science of Cause and Effect*. Basic Books, Inc., USA, 1st edition.
- Pearl, J. and Paz, A. (1986). Graphoids: Graph-based logic for reasoning about relevance relations or when would x tell you more about y if you already know z? In *ECAI*.
- Reinis, Z., Pokorný, J., Basika, V., Tiserova, J., Gorican, K., Horakova, D., Stuchlikova, E., Havranek, T., and Hrabovsky, F. (1981). Prognostic significance of the risk profile in the prevention of coronary heart disease. *Bratisl Lek Listy*, 76:137–150.
- Runge, J., Bathiany, S., Bollt, E., Camps-Valls, G., Coumou, D., Deyle, E., Glymour, C., Kretschmer, M., Mahecha, M. D., Muñoz-Marí, J., et al. (2019). Inferring causation from time series in earth system sciences. *Nature communications*, 10(1):1–13.

- Sadeghi, K. and Lauritzen, S. (2014). Markov Properties for Mixed graphs. *Bernoulli*, 20(2):676 – 696.
- Schölkopf, B. (2019). Causality for machine learning. *arXiv preprint arXiv:1911.10500*.
- Scholz, F. W. and Stephens, M. A. (1987). K-sample anderson-darling tests. *Journal of the American Statistical Association*, 82(399):918–924.
- Shah, R. D. and Peters, J. (2020). The hardness of conditional independence testing and the generalised covariance measure. *The Annals of Statistics*, 48(3).
- Silander, T. and Leong, T.-Y. (2013). A dynamic programming algorithm for learning chain event graphs. In *International Conference on Discovery Science*, pages 201–216. Springer.
- Simpson, E. H. (1951). The interpretation of interaction in contingency tables. *Journal of the Royal Statistical Society: Series B (Methodological)*, 13(2):238–241.
- Solus, L., Wang, Y., and Uhler, C. (2021). Consistency Guarantees for Greedy Permutation-Based Causal Inference Algorithms. *Biometrika*. asaa104.
- Spirtes, P. and Glymour, C. (1991). An algorithm for fast recovery of sparse causal graphs. *Social Science Computer Review*, 9(1):62–72.
- Spirtes, P., Glymour, C., and Scheines, R. (2000). *Causation, Prediction, and Search*. MIT press, 2nd edition.
- Tarjan, R. E. (1976). Edge-disjoint spanning trees and depth-first search. *Acta Informatica*, 6(2):171–185.
- Thwaites, P., Smith, J. Q., and Riccomagno, E. (2010). Causal analysis with chain event graphs. *Artificial Intelligence*, 174(12):889–909.
- Tsamardinos, I., Brown, L. E., and Aliferis, C. F. (2006). The max-min hill-climbing bayesian network structure learning algorithm. *Machine learning*, 65(1):31–78.
- Tu, R., Zhang, C., Ackermann, P., Mohan, K., Kjellström, H., and Zhang, K. (2019). Causal discovery in the presence of missing data. In Chaudhuri, K. and Sugiyama, M., 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.
- Uhler, C., Raskutti, G., Bühlmann, P., and Yu, B. (2013). Geometry of the Faithfulness Assumption in Causal inference. *The Annals of Statistics*, 41(2):436 – 463.

- Verma, T. and Pearl, J. (1990). *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.
- Verma, T. S. and Pearl, J. (2013). On the equivalence of causal models. *CoRR*, abs/1304.1108.
- Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S. J., Brett, M., Wilson, J., Millman, K. J., Mayorov, N., Nelson, A. R. J., Jones, E., Kern, R., Larson, E., Carey, C. J., Polat, İ., Feng, Y., Moore, E. W., VanderPlas, J., Laxalde, D., Perktold, J., Cimrman, R., Henriksen, I., Quintero, E. A., Harris, C. R., Archibald, A. M., Ribeiro, A. H., Pedregosa, F., van Mulbregt, P., and SciPy 1.0 Contributors (2020). SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272.
- Wagner, C. H. (1982). Simpson’s paradox in real life. *The American Statistician*, 36(1):46–48.
- Yang, K., Katcoff, A., and Uhler, C. (2018). Characterizing and learning equivalence classes of causal DAGs under interventions. In Dy, J. and Krause, A., editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5541–5550. PMLR.
- Zheng, X., Aragam, B., Ravikumar, P. K., and Xing, E. P. (2018). Dags with no tears: Continuous optimization for structure learning. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.

Appendix A : Implementation remarks

All of the code used in this work is available at github.com/mnazaal/masters-thesis. We make extensive use of `numpy` ¹, `scipy` ² (for statistical testing), `networkx` ³ (for graphs), `pgmpy` ⁴ (for initial DAG learning) in Python for the work on CStrees and the `stagedtrees` package ⁵ in R to learn the staged trees.

When using the Epps-Singleton test implemented in `scipy`, one must add an epsilon to the difference between the upper and lower quartile within the implementation since the test performs a computation involving division by this value. For the KL divergence we also add an epsilon to the estimated probability density in order to avoid dividing by 0 which can happen if there are no samples taking a specific value.

Our usage indicated that the Hill climbing algorithm is highly susceptible to local minima when optimizing for the BIC score.

We make use of `pytest` ⁶ to run unit tests. This thesis is written in `org-mode` with the intention for users to be able to execute the code during the compilation of the document itself to ensure full reproducibility.

¹ Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M. H., Brett, M., Haldane, A., del Río, J. F., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K., Reddy, T., Weckesser, W., Abbasi, H., Gohlke, C., and Oliphant, T. E. (2020). Array programming with NumPy. *Nature*, 585(7825):357–362

² Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S. J., Brett, M., Wilson, J., Millman, K. J., Mayorov, N., Nelson, A. R. J., Jones, E., Kern, R., Larson, E., Carey, C. J., Polat, İ., Feng, Y., Moore, E. W., VanderPlas, J., Laxalde, D., Perktold, J., Cimrman, R., Henriksen, I., Quintero, E. A., Harris, C. R., Archibald, A. M., Ribeiro, A. H., Pedregosa, F., van Mulbregt, P., and SciPy 1.0 Contributors (2020). SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272

³ Hagberg, A. A., Schult, D. A., and Swart, P. J. (2008). Exploring network structure, dynamics, and function using `networkx`. In Varoquaux, G., Vaught, T., and Millman, J., editors, *Proceedings of the 7th Python in Science Conference*, pages 11 – 15, Pasadena, CA USA

⁴ Ankan, A. and Panda, A. (2015). `pgmpy`: Probabilistic graphical models using python. In *Proceedings of the 14th Python in Science Conference (SCIPY 2015)*. Citeseer

⁵ Carli, F., Leonelli, M., Riccomagno, E., and Varando, G. (2020). The `r` package `stagedtrees` for structural learning of stratified staged trees. *arXiv preprint arXiv:2004.06459*

⁶ Krekel, H., Oliveira, B., Pfannschmidt, R., Bruynooghe, F., Laughner, B., and Bruhin, F. (2004). `pytest x.y`

Appendix B : Extra figures

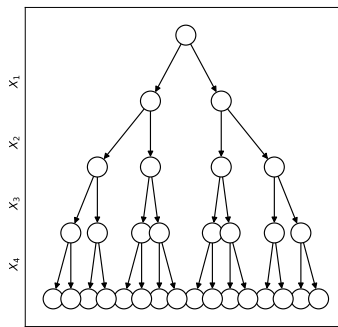
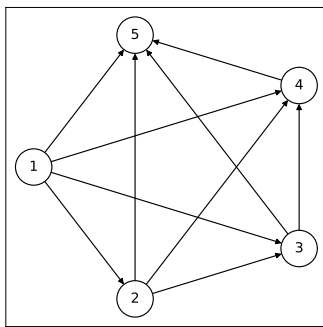


Figure 8.1: Generating the CStree for a fully connected DAG using the ordering 1234 with Algorithm 1

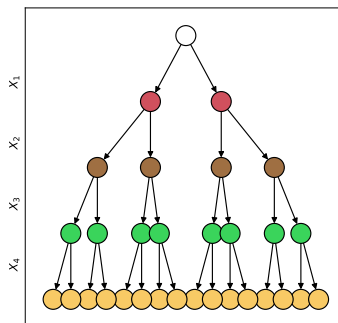
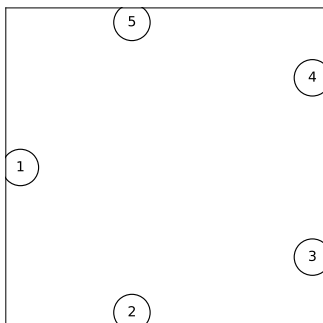


Figure 8.2: Generating the CStree for the empty DAG using the ordering 1234 with Algorithm 1

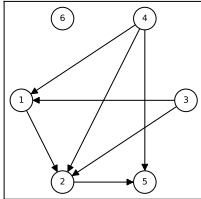
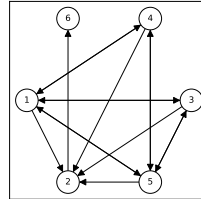
PC**Hill climbing**

Figure 8.3: CPDAG representing the MEC of the DAG learnt from the coronary dataset. The PC algorithm suggested 12 causal orderings, whilst the Hill climbing algorithm suggested 16.

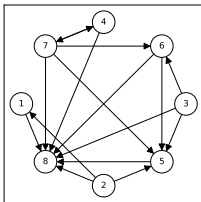
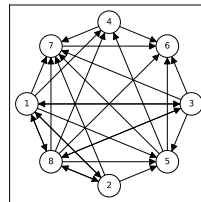
PC**Hill climbing**

Figure 8.4: CPDAG representing the MEC of the DAG learnt from the mice protein expression dataset. The PC algorithm suggested 196 causal orderings, whilst the Hill climbing algorithm suggested 12.

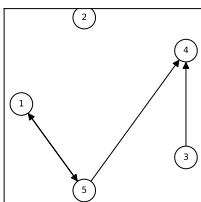
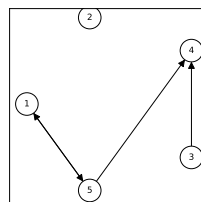
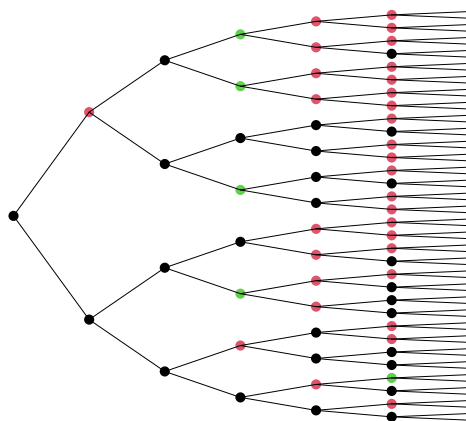
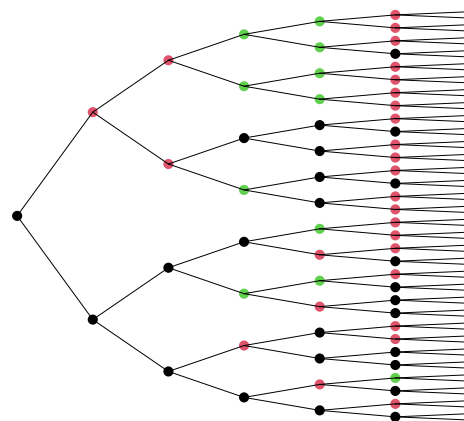
PC**Hill climbing**

Figure 8.5: CPDAG representing the MEC of the DAG learnt from the Vitamin D dataset. Both algorithms gave the same CPDAG and suggested 40 possible causal orderings.

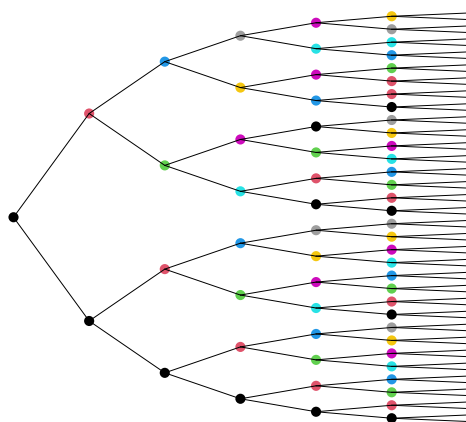
Hill Climbing



Backwards hill climbing



Backwards joining



Hierarchical clustering

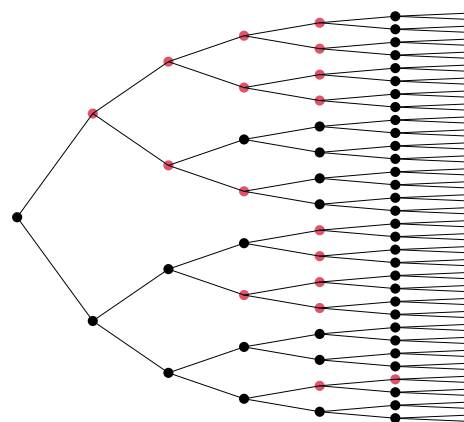
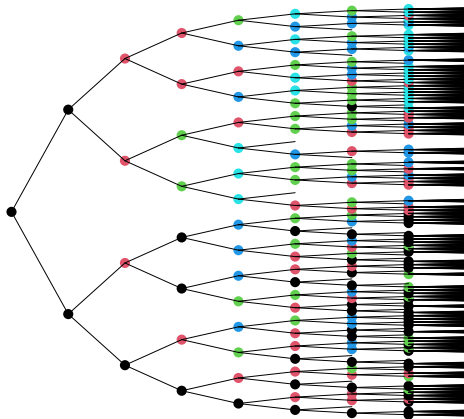
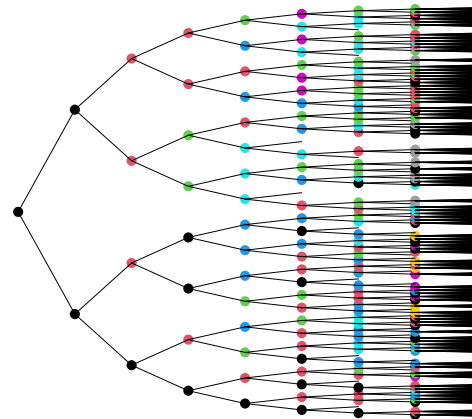


Figure 8.6: BIC Optimal staged trees learnt from the coronary dataset using the ordering 534126.

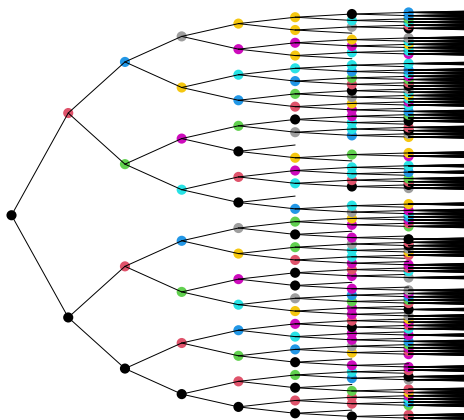
Hill Climbing



Backwards hill climbing



Backwards joining



Hierarchical clustering

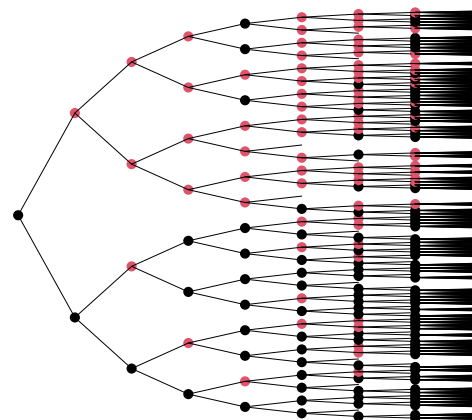
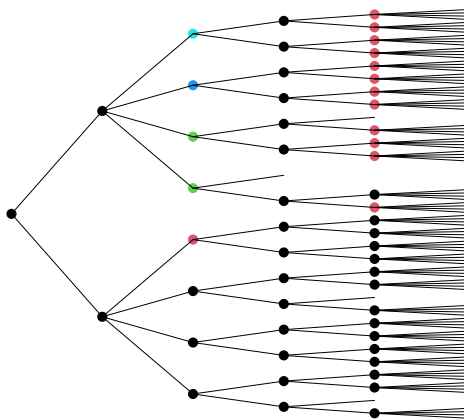
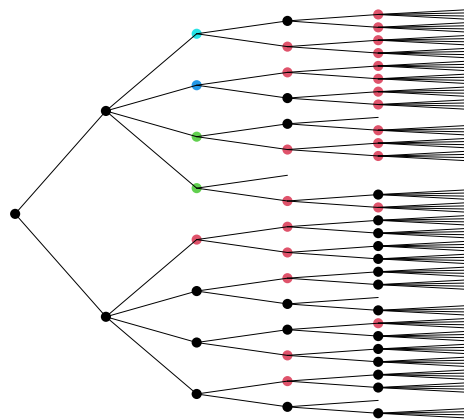


Figure 8.7: BIC Optimal staged trees learnt from the mice protein expression dataset using the ordering 32765148.

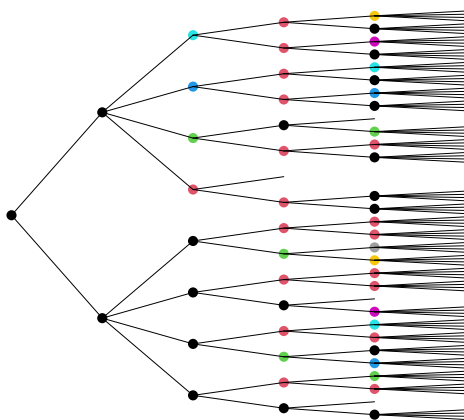
Hill Climbing



Backwards hill climbing



Backwards joining



Hierarchical clustering

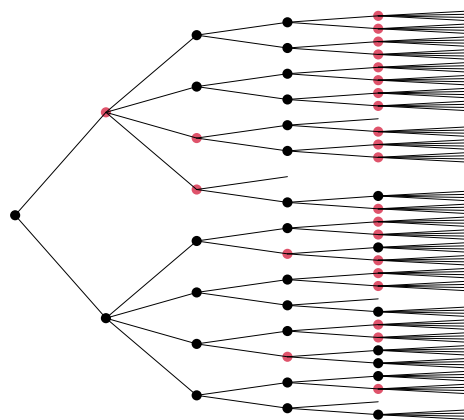
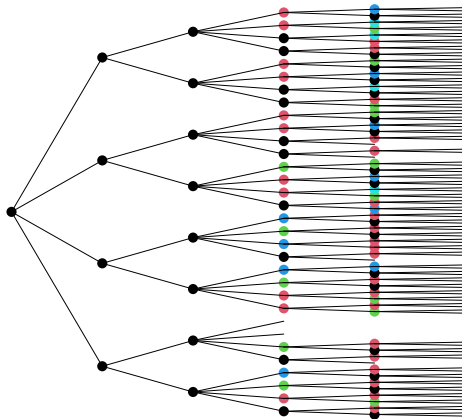
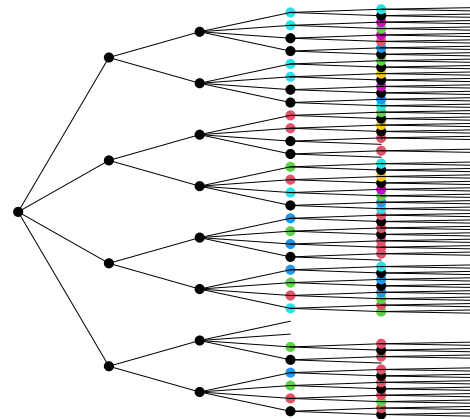


Figure 8.8: BIC Optimal staged trees learnt from the Vitamin D dataset using the ordering 53421.

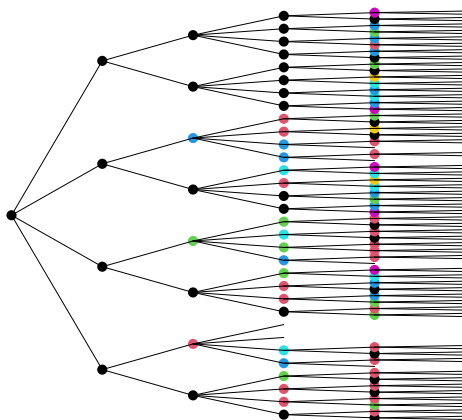
Hill Climbing



Backwards hill climbing



Backwards joining



Hierarchical clustering

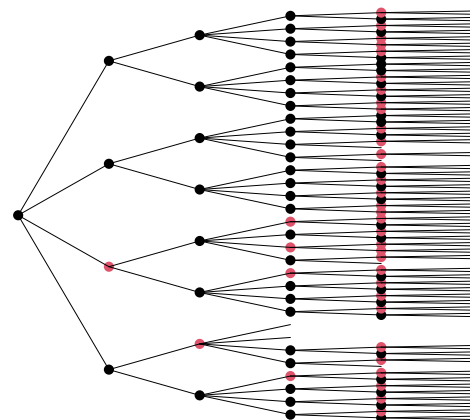


Figure 8.9: BIC Optimal staged trees learnt from the Vitamin D dataset using the ordering 32145.