

A Weaker Faithfulness Assumption based on Triple Interactions

Alexander Marx
MPI-INF and CISPA*
Saarland University
Saarbrücken, Germany

Arthur Gretton
Gatsby Unit
University College London
London, United Kingdom

Joris M. Mooij
Korteweg-de Vries Institute
University of Amsterdam
Amsterdam, The Netherlands

Abstract

One of the core assumptions in causal discovery is the faithfulness assumption—i.e. assuming that independencies found in the data are due to separations in the true causal graph. This assumption can, however, be violated in many ways, including xor connections, deterministic functions or cancelling paths. In this work, we propose a weaker assumption that we call 2-adjacency faithfulness. In contrast to adjacency faithfulness, which assumes that there is no conditional independence between each pair of variables that are connected in the causal graph, we only require no conditional independence between a node and a subset of its Markov blanket that can contain up to two nodes. Equivalently, we adapt orientation faithfulness to this setting. We further propose a sound orientation rule for causal discovery that applies under weaker assumptions. As a proof of concept, we derive a modified Grow and Shrink algorithm that recovers the Markov blanket of a target node and prove its correctness under strictly weaker assumptions than the standard faithfulness assumption.

1 INTRODUCTION

In this work, we focus on causal discovery from observational data, where we are given a sample from the joint distribution P of the observed variables and try to infer the true causal graph G between them. Two standard assumptions in this field are the causal Markov condition and the faithfulness assumption (Spirtes et al., 2000). While the causal Markov condition assumes that all separations in the true causal graph G imply independencies in P , the faithfulness assumption is its counterpart.

*Max Planck Institute for Informatics and CISPA Helmholtz Center for Information Security

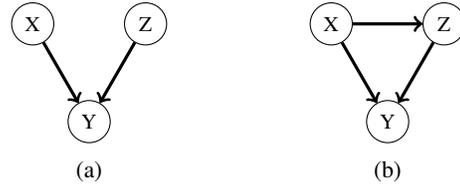


Figure 1: Failures of adjacency faithfulness: Assume in graph (a) X, Z are fair independent coins and $Y := (X \oplus Z) \oplus E$, where \oplus is the xor operator and E is a biased coin denoting a noise term. Then X is independent of Y (denoted as $X \perp_P Y$) and $Z \perp_P Y$. Graph (b) could correspond with a linear model where both directed paths from X to Y cancel s.t. $X \perp_P Y$, but $X \not\perp_P Z$ and $Z \not\perp_P Y$.

That is, all independencies found in P are due to separations in G . Although both assumptions have great merit for causal discovery algorithms, especially the faithfulness assumption has been criticized in the past (Anderesen, 2013; Zhang and Spirtes, 2016).

Despite it was proven that faithfulness violations in causally sufficient linear-Gaussian and discrete acyclic systems occur with Lebesgue measure zero (Meek, 1995b), it has also been shown that given a finite sample, empirical faithfulness violations do appear surprisingly often (Uhler et al., 2013). Even on population level, there exist simple generating mechanisms, as shown in Figure 1, that violate the faithfulness assumption. For instance, two independent random variables X and Z , that can be modelled by fair coins, together cause Y through a noisy xor relation. As a consequence, all three variables are marginally independent. Following the faithfulness assumption, there should be no edges connecting X, Y and Z in the causal graph—however, there are.

Faithfulness violations like the above have been intensively studied in the past (Ramsey et al., 2006; Zhang and Spirtes, 2008; Spirtes and Zhang, 2014) and several weaker assumptions such as adjacency

faithfulness (Spirtes et al., 2000), P-minimality (Pearl, 2009), SGS-minimality (Spirtes et al., 2000) and frugality (Forster et al., 2017), which we review in Section 3.3, have been proposed. Although faithfulness violations induced by xor-type relations—i.e. both parents are marginally independent of the child node—can be detected by most of the above approaches, they do not analyze under which conditions the DAG structure can be recovered, once such violations have been detected.

In this work, we propose a new assumption that we call *2-adjacency faithfulness*, which allows us to both detect such faithfulness violations and partially infer the underlying DAG structure under certain conditions. We start by explaining the standard concepts and notation in Section 2 and review failures of adjacency faithfulness as well as related work in Section 3. Then, we study the causal structure of xor-type connections in Section 4 and propose 2-adjacency faithfulness in Section 5. To partially infer causal DAGs that may contain such generating mechanisms, we introduce a sound orientation rule, in Section 6. Further, we show under which assumptions on the distribution this rule is applicable—which we formalize as the 2-orientation faithfulness assumption—and analyze its failure cases. As a proof of concept, we introduce a modification of the Grow and Shrink (GS) algorithm (Margaritis and Thrun, 2000) in Section 7 and show it correctly identifies the Markov blanket of a target node under strictly weaker assumptions than faithfulness. In addition, we give some intuition on how to extend well known causal discovery algorithms based on our new assumptions.

2 DAGS AND INDEPENDENCE

In this section, we define our notation and provide definitions for separations on graphs and independence.

2.1 Causal Graphs

A causal *directed acyclic graph* (DAG) G over a set of random variables V with joint distribution P is defined such that each pair of nodes that is adjacent in G is causally related. For simplicity, we will use the random variables V to also refer to the nodes of the graph. A directed edge $X \rightarrow Y$ in G between two nodes representing the random variables $X, Y \in V$ indicates that X is a *direct cause* or *parent* of Y and that Y is a *direct effect* or *child* of X . Accordingly, we denote the set of all parents of $X \in V$ with $\text{Pa}(X)$, the set of all children with $\text{Ch}(X)$ and the set of parents and children with $\text{PC}(X) := \text{Pa}(X) \cup \text{Ch}(X)$. Further, we write $\text{An}(X)$ for the set of *ancestors* and $\text{De}(X)$ for the set of all *descendants* of X , where X is an ancestor and descendant of

itself. Respectively, we refer to the *non-descendants* of X as $\text{Nd}(X) := V \setminus \text{De}(X)$. Last, the *Markov blanket* of a variable X is defined as $\text{MB}(X) := \text{PC}(X) \cup \text{Sp}(X)$, where $\text{Sp}(X)$ are the spouses of X , that is, nodes that share a child node with X . Importantly, X is d -separated of any other node in the graph given its Markov blanket and $\text{MB}(X)$ is the smallest such set.

DAGs are used to represent causal graphs under the assumption of acyclicity, no selection bias, and *causal sufficiency*, that is, it is assumed that no two variables $X, Y \in V$ are caused by a confounder Z which is not in the set of observed variables V . This is also the setup on which we focus in this paper—i.e. assuming that all relevant variables are observed, that there are no causal cycles and that there has been no conditioning on selection variables. Further, as a short form to summarize a model as defined above, we write $\mathcal{M} = (G, V, P)$.

2.2 Independence and Separation

In the following, we define conditional independence in a probability distribution and d -separation in a graph.

Given three sets of probabilistic random variables $X, Y, Z \subseteq V$, where P is the joint distribution over V , we denote that X is *probabilistically independent* of Y given Z in P as $X \perp\!\!\!\perp_P Y \mid Z$.

d-separation (Pearl, 2009) is defined in terms of paths. A *path* p between X and Y , denoted $p = \langle X, \dots, Y \rangle$, is a sequence of distinct nodes X_1, \dots, X_n such that X_i is adjacent to X_{i+1} for $i = 1, \dots, n - 1$, $X_1 = X$ and $X_n = Y$. Further, we call a node C a *collider* on a path $\langle \dots, X, C, Y, \dots \rangle$, where C is adjacent to both X and Y , if two arrowheads point to it, that is $X \rightarrow C \leftarrow Y$.

Definition 1 (*d*-Separation) A path between two vertices X, Y in a DAG is *d*-connecting given a set Z , if

1. every non-collider on the path is not in Z , and
2. every collider on the path is an ancestor of Z .

If there is no path *d*-connecting X and Y given Z , then X and Y are *d*-separated given Z . Sets X and Y are *d*-separated given Z , if for every pair X, Y , with $X \in X$ and $Y \in Y$, X and Y are *d*-separated given Z .

As shorthand notation for separations on a DAG G , we write $X \perp\!\!\!\perp_G Y \mid Z$ if X is *d*-separated from Y given Z . Following this notation, we state the graphoid axioms (Dawid, 1979; Spohn, 1980; Geiger et al., 1990).

Definition 2 (Graphoid Axioms) Let $\mathcal{M} = (G, V, P)$, with $W, X, Y, Z \subseteq V$. The (semi-)graphoid axioms are the following rules ($\perp\!\!\!\perp$ denotes $\perp\!\!\!\perp_P$ and $\perp\!\!\!\perp_G$)

1. *Symmetry*: $X \perp\!\!\!\perp Y \mid Z \Rightarrow Y \perp\!\!\!\perp X \mid Z$.

2. *Decomposition:* $X \perp\!\!\!\perp Y \cup W \mid Z \Rightarrow X \perp\!\!\!\perp Y \mid Z$.
3. *Weak Union:* $X \perp\!\!\!\perp Y \cup W \mid Z \Rightarrow X \perp\!\!\!\perp Y \mid W \cup Z$.
4. *Contraction:* $(X \perp\!\!\!\perp Y \mid W \cup Z) \wedge (X \perp\!\!\!\perp W \mid Z) \Rightarrow X \perp\!\!\!\perp Y \cup W \mid Z$.

For separations only on the graph, the graphoid axioms include two additional rules (only for \perp_G).

5. *Intersection:* $(X \perp\!\!\!\perp Y \mid W \cup Z) \wedge (X \perp\!\!\!\perp W \mid Y \cup Z) \Rightarrow X \perp\!\!\!\perp Y \cup W \mid Z$, for any pairwise disjoint subsets $W, X, Y, Z \subseteq V$.
6. *Composition:* $(X \perp\!\!\!\perp Y \mid Z) \wedge (X \perp\!\!\!\perp W \mid Z) \Rightarrow X \perp\!\!\!\perp Y \cup W \mid Z$.

As an illustration why certain rules only hold for graphs and not generally for probability distributions, consider rule (6) and Figure 1 (a) again. From the distribution induced by the xor, we find that $Y \perp\!\!\!\perp_P X$ and $Y \perp\!\!\!\perp_P Z$ but we cannot conclude that $Y \perp\!\!\!\perp_P \{X, Z\}$. If, however, in a graph Y is d -separated from X and from Z then Y is d -separated from the set $\{X, Z\}$.

We round up this section by defining the causal Markov condition (Spirtes et al., 2000) (CMC) for DAGs.¹

Definition 3 (Causal Markov Condition) Given the triple $\mathcal{M} = (G, V, P)$, the causal Markov condition holds, if every d -separation imposed by G implies an independence in P .

The causal Markov condition is one of the most essential assumptions for causal discovery algorithms. On the other hand, assumptions about what properties of the graph can be inferred based on the given distribution have been weakened over time (Ramsey et al., 2006; Zhang and Spirtes, 2008; Forster et al., 2017). Most commonly known is the faithfulness assumption.

3 ADJACENCY FAITHFULNESS AND WHEN IT IS VIOLATED

To lay out the problem, we first explain faithfulness and adjacency faithfulness, then examine when those could fail and give a summary about the most relevant related approaches that use weaker assumptions.

The faithfulness assumption is one of the core assumptions made by most causal discovery algorithms (Spirtes et al., 2000) and it can be seen as the inverse assumption to CMC—i.e. assuming that all independencies found in P imply a d -separation in the causal graph. Adjacency faithfulness is a slightly weaker assumption.

¹CMC is also often referred to as the global Markov condition, which is equivalent to the local Markov condition for all dependency models that obey the semi-graphoid axioms (Pearl and Verma, 1987).

Definition 4 (Adjacency Faithfulness) Given $\mathcal{M} = (G, V, P)$, if $X, Y \in V$ are adjacent in G , then they are probabilistically dependent given all $S \subseteq V \setminus \{X, Y\}$.

Alternatively, we could turn this definition around by stating that if we find a conditional independence in P , then we assume that there is no edge in the corresponding graph. Assuming adjacency faithfulness ensures that we recover the correct skeleton graph (i.e. the undirected graph). Correct detection of the skeleton together with the correct identification of all collider structures ensures that the detected graph is in the Markov equivalence class of the true graph (Verma and Pearl, 1991). The missing ingredient—i.e. the correct detection of all collider structures—is ensured by additionally assuming orientation faithfulness (Zhang and Spirtes, 2008).

Definition 5 (Orientation-Faithfulness) Given $\mathcal{M} = (G, V, P)$. Let the path $\langle X, Y, Z \rangle$ be unshielded² in G .

1. If $X \rightarrow Y \leftarrow Z$, then X and Z are dependent given any subset in $V \setminus \{X, Z\}$ that contains Y ; otherwise
2. X and Z are dependent conditional on any subset of $V \setminus \{X, Z\}$ that does not contain Y .

The bottleneck here is the adjacency faithfulness assumption, as many causal discovery algorithms such as PC (Spirtes et al., 2000) or GES (Chickering, 2002) rely on finding adjacent nodes either by checking for marginal dependencies or adding single edges based on adjacency faithfulness and CMC. If one is willing to assume that those assumptions hold, then any violation of orientation faithfulness can be detected as shown by Zhang and Spirtes (2008). However, adjacency faithfulness can be violated in many ways, e.g. by xor-type connections, path cancellations, or deterministic relations. We briefly explain the first two below, as they are relevant for the remainder. For deterministic relations and finite sample failures, we refer to Lemeire et al. (2012).

3.1 Xor-Type Relations

In this work, we focus on xor-type relations. That is, given a triple of nodes $X, Y, Z \in V$ such that $X \rightarrow Y \leftarrow Z$, where at least one of the causal edges cannot be detected by marginal dependence, but only by looking at the joint distribution over X, Y and Z . The key here is that either parent of Y might not be dependent on Y , but only by considering both parents, we can detect the dependence. To illustrate this, consider the following example where we describe a noisy xor with an unobserved noise variable that is modelled with a biased coin as it is common for binary causal structures (Inazumi et al., 2011).

²For an unshielded path $\langle X, Y, Z \rangle$, X is adjacent to Y and Y is adjacent to Z , but X is not adjacent to Z .

Example 1 Let $M = (G, V, P)$ be a causal model. Given variables $X, Y, Z \in V$ such that $X \rightarrow Y \leftarrow Z$ in G and there is no edge connecting X and Z , as in Figure 1(a), where X, Z are fair independent coins. Their common effect Y is generated as $Y := (X \oplus Z) \oplus E$, where \oplus denotes xor—i.e. $X \oplus Z := (X + Z) \bmod 2$ —and E is a biased coin with $P(E = 1) = p$, where $0 \leq p < \frac{1}{2}$ and $E \perp\!\!\!\perp_P \{X, Z\}$. Hence, $X \not\perp\!\!\!\perp_G Y$, $Z \not\perp\!\!\!\perp_G Y$, however, due to the xor, we have that $X \perp\!\!\!\perp_P Y$ and $Z \perp\!\!\!\perp_P Y$. Both edges violate adjacency faithfulness. If we were to check the joint distribution, we can find that $Y \not\perp\!\!\!\perp_P \{X, Z\}$, or $X \not\perp\!\!\!\perp_P Z \mid Y$, since we get that $P(X = 1, Z = 1, Y = 1) = \frac{p}{4}$, where $P(X = 1, Z = 1) \cdot P(Y = 1) = \frac{1}{4} \cdot \frac{1}{2} = \frac{1}{8}$. Those terms being equal would only hold if $p = \frac{1}{2}$, which we excluded by assumption.

Similar examples, where the marginal dependencies might be hard to detect can be found for continuous data (Sejdicinovic et al., 2013)—e.g. if X, Z are normally distributed with mean zero and variance one, and $Y := \text{sign}(XZ) \cdot E$, with exponentially distributed noise $E \sim \text{Exp}(\frac{1}{\sqrt{2}})$ (see Figure 2).

3.2 Cancelling Paths

A minimal example of cancelling paths was given by Hesslow (1976) and is illustrated with the causal graph shown in Figure 1(b). In Hasslow’s example taking birth control pills (X) can influence the risk of getting thrombosis (Y) via two paths. It has a direct effect and also taking the pills reduces the chance of pregnancy (Z), which itself is a cause of thrombosis. However, the causal effects induced by those paths cancel such that $X \perp\!\!\!\perp_P Y$ even though $X \not\perp\!\!\!\perp_G Y$. This failure of faithfulness was shown to be undetectable since X will be dependent on Y given Z and hence the graph $X \rightarrow Z \leftarrow Y$ is also a valid graph for those independencies—i.e. Markov equivalent (Zhang and Spirtes, 2008). There exist cancelling paths that consist of more than three variables, which are detectable, e.g. if Z is not adjacent to Y , but there is a path $Z \rightarrow W \rightarrow Y$ (Zhang and Spirtes, 2008).

3.3 Weaker Assumptions

In the following, we discuss different approaches on how to relax the faithfulness assumption.

Two well-studied assumptions are P-minimality (Pearl, 2009) and SGS-minimality (Spirtes et al., 2000). While the former states that from all DAGs that satisfy the causal Markov condition w.r.t. P , the DAG that entails most conditional independence statements is preferred. The latter assumes that no proper subgraph of the true DAG fulfils the causal Markov condition w.r.t. to P . From both assumptions, SGS-minimality is the weaker

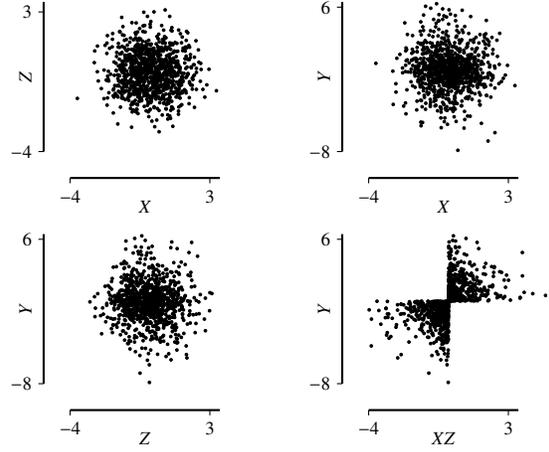


Figure 2: Sample data for the collider graph $X \rightarrow Y \leftarrow Z$, where $X, Z \sim N(0, 1)$ are iid and $Y := \text{sign}(XZ) \cdot E$, with $E \sim \text{Exp}(\frac{1}{\sqrt{2}})$. The dependence is only detectable by considering all three variables jointly.

assumption (Zhang, 2013). In a different line of research, it was shown that SGS-minimality suffices for causal discovery approaches based on the additive noise assumption (Peters et al., 2014).

A more recent approach by Forster et al. (2017) introduces the concept of *frugality*, which is a stronger assumption than both minimality assumptions. The authors define a DAG G to be more frugal than G' , if G contains fewer edges than G' . A maximally frugal DAG uses only as many edges as are necessary to satisfy the causal Markov condition. To determine maximally frugal graphs, one has to consider all causal orderings of the variables, which is rather costly, but can be solved using permutation algorithms (Raskutti and Uhler, 2018). Another approach to discover causal graphs based on frugality, or any of the above assumptions is based on boolean satisfiability (SAT) solvers (Zhalama et al., 2017).

In this paper, we introduce the 2-adjacency faithfulness assumption, which allows us to find xor-type relations, some faithfulness violations induced by cancelling paths and all relations that can be detected by assuming adjacency faithfulness. We conjecture that 2-adjacency faithfulness is a slightly stronger assumption than frugality since frugality considers all permutations and not only triples (Forster et al., 2017). However, this can also be an advantage, since we only need to check all triples to detect 2-associations, whereas the most frugal graph can only be determined by considering all permutations. In addition, we extend existing work by providing a sound orientation rule that can be used to infer the edges within a 2-association, if they appear in a larger graph.

In the next section, we discuss xor-type relations, which are a generalization of Example 1 and can be described as 2-associations. We use those structures as an example to illustrate one of the main properties of 2-associations, that we describe in Theorem 1.

4 UNFAITHFUL TRIPLES

We first define what we call an unfaithful triple³ and its properties, and then argue why such a triple a) violates adjacency faithfulness and b) even if detected, the underlying DAG structure cannot be uniquely determined without further information.

Definition 6 (Unfaithful Triple) Given $\mathcal{M} = (G, \mathbf{V}, P)$ and three distinct nodes $X, Y, Z \in \mathbf{V}$: if X, Y and Z are marginally independent but not mutually independent in P , we call $\{X, Y, Z\}$ an unfaithful triple w.r.t. P .⁴ If further for each distinct pair of nodes $A, B \in \{X, Y, Z\}$:

$$\forall S \subseteq \mathbf{V} \setminus \{X, Y, Z\} : A \not\perp_P B \mid S \cup \{X, Y, Z\} \setminus \{A, B\},$$

we call $\{X, Y, Z\}$ a minimal unfaithful triple.

The first example for such a triple for three binary random variables was given by Bernstein (1927), which is equivalent to our noisy xor example. The minimality condition ensures that the three nodes are connected by a path of length two, as we will show below. This concept is also illustrated in Figure 3.

We start by showing that if three random variables $\{X, Y, Z\}$ are marginally independent, finding a dependence between all three variables, e.g. $X \not\perp_P \{Y, Z\}$, implies that also $Y \not\perp_P \{X, Z\}$ and $Z \not\perp_P \{X, Y\}$.

Lemma 1 Given $\mathcal{M} = (G, \mathbf{V}, P)$, let $\{X, Y, Z\} \subseteq \mathbf{V}$ form an unfaithful triple in P , then $X \not\perp_P \{Y, Z\}$, $Y \not\perp_P \{X, Z\}$ and $Z \not\perp_P \{X, Y\}$, which in addition implies that $X \not\perp_P Y \mid Z$, $X \not\perp_P Z \mid Y$ and $Y \not\perp_P Z \mid X$.

Proof: Assume that w.l.o.g. $X \not\perp_P \{Y, Z\}$ is violated. By weak union, we get $X \perp_P Y \mid Z$ which is equivalent to $Y \perp_P X \mid Z$, using symmetry. We know that $Y \perp_P Z$. By contraction, we get that $Y \perp_P \{X, Z\}$. Similarly, we conclude that $Z \perp_P \{X, Y\}$. Altogether, this implies that X, Y, Z would be independent, which is a contradiction.

Each pair of joint dependence and marginal independence, e.g. $X \not\perp_P \{Y, Z\}$ and $X \perp_P Z$, implies a conditional dependence, e.g. $X \not\perp_P Y \mid Z$, by contraction. \square

³Ramsey et al. (2006) used the term unfaithful triple for the non-detectable faithfulness violation explained in Section 3.2.

⁴Not mutually independent implies that $X \not\perp_P \{Y, Z\}$, $Y \not\perp_P \{X, Z\}$ or $Z \not\perp_P \{X, Y\}$.

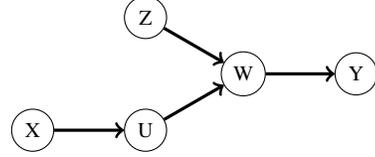


Figure 3: Assume that $\{X, Y, Z\}$ form an unfaithful triple. Since X is d -separated from Y given U and Z , they do not form a minimal unfaithful triple. Neither do $\{U, Y, Z\}$, since Y can be d -separated from U given $\{W, Z\}$. Thus, only $\{U, W, Z\}$ can be a minimal unfaithful triple.

Consider Example 1. Since X, Y, Z form an unfaithful triple, we can infer from Lemma 1 that each pair is conditionally dependent given the third node. As there are no other nodes in the graph, X, Y, Z must form a minimal unfaithful triple.

Next, we show that (minimal) unfaithful triples must be connected in the causal graph.

Lemma 2 Given $\mathcal{M} = (G, \mathbf{V}, P)$, let $\{X, Y, Z\} \subseteq \mathbf{V}$ form an unfaithful triple in P . If CMC holds, each node in the triple is d -connected to at least one other node in the triple by a path in G .

Proof: Assume w.l.o.g. that X is d -separated from Y and Z in G —i.e. $X \perp_G Y$ and $X \perp_G Z$. By applying the composition axiom, we get that $X \perp_G \{Y, Z\}$. If we apply the causal Markov condition, we get that $X \perp_P \{Y, Z\}$, which is a contradiction to our assumption. \square

Further, we show that a minimal unfaithful triple has to contain a collider on a path of length two that connects all three nodes in the triple, e.g. $X \rightarrow Y \leftarrow Z$. To do that, we first show a more general statement.

Theorem 1 Given $\mathcal{M} = (G, \mathbf{V}, P)$ with three distinct nodes $X, Y, Z \in \mathbf{V}$ and assume that CMC holds. If $\forall S \subseteq \mathbf{V} \setminus \{X, Y, Z\}$ it holds that $X \not\perp_P Y \mid Z \cup S$, $X \not\perp_P Z \mid Y \cup S$ and $Y \not\perp_P Z \mid X \cup S$, then one of the three nodes is a collider on a path of length two between the two other nodes, e.g. $X \rightarrow Y \leftarrow Z$ in G .

Proof: There must be (at least) one node in $\{X, Y, Z\}$ that is not an ancestor of any of the other nodes, say $Z \notin \text{An}(X)$ and $Z \notin \text{An}(Y)$, because of acyclicity. In other words, $X \notin \text{De}(Z)$ and $Y \notin \text{De}(Z)$. The local Markov property states that $Z \perp_G \text{Nd}(Z) \mid \text{Pa}(Z)$ and hence in particular $Z \perp_G \{X, Y\} \mid \text{Pa}(Z)$. Further, if $|\text{Pa}(Z) \cap \{X, Y\}| < 2$, we get a contradiction with the assumed conditional dependences. Hence $\{X, Y\} \subseteq \text{Pa}(Z)$ and $X \rightarrow Z \leftarrow Y$ is in G . \square

The theorem only states that there exists a collider, e.g. $X \rightarrow Y \leftarrow Z$, but not whether this path is shielded or not.

Since we do not assume any marginal dependence or independence in Theorem 1, we can derive that the same statement holds for a minimal unfaithful triple. Notice that for a minimal unfaithful triple each pair of nodes is marginally independent, which implies that there is no way to decide which of the three possible collider structures corresponds with the causal graph in the absence of further information.

Knowing that a minimal unfaithful triple has to contain a collider in G , it is obvious that such a structure violates adjacency faithfulness, as none of the edges is represented by a marginal dependence in P . The key point is that we can detect such interactions by taking multiple parents into account. In the following, we define a weaker assumption that allows us to detect and infer causal graphs that contain such faithfulness violations.

5 2-ADJACENCY FAITHFULNESS

To define our new assumption, we first need to define associations between a single node and a set of nodes.

Definition 7 (k -Association) *Let P be the joint distribution of a set of observed random variables V .*

1. *Given distinct $X, Y \in V$, we say that X is 1-associated to Y , if $\forall S \subseteq V \setminus \{X, Y\} : X \not\perp_P Y \mid S$.*
2. *Given distinct $X, Y_1, Y_2 \in V$, X is 2-associated to $\{Y_1, Y_2\}$ if $\forall S \subseteq V \setminus \{X, Y_1, Y_2\}$ it holds that*
 - i) $X \not\perp_P Y_1 \mid S \cup Y_2$,
 - ii) $X \not\perp_P Y_2 \mid S \cup Y_1$ and
 - iii) $Y_1 \not\perp_P Y_2 \mid S \cup X$.

We call X strictly 2-associated to $\{Y_1, Y_2\}$, if X is 2-associated to $\{Y_1, Y_2\}$ and not 1-associated to Y_1 or Y_2 .

In other words, k -associations relate to two types of dependencies: certain conditional dependencies between pairs of variables (1-associations) and between triples (2-associations). For readability, we use a shorthand notation and write $X \perp_2 \{Y, Z\}$ if X is 2-associated to Y and Z resp. $X \perp_1 Y$ if X is 1-associated to Y . If we refer to a set Y that contains at most two elements and we want to express that X is either 1- or 2-associated to this set, we write $X \perp_{\leq 2} Y$. We denote a strict 2-association by “ $\overset{s}{\perp}_2$ ”.

Pairwise dependencies can occur for example in a simple chain $X \rightarrow Y \rightarrow Z$, where no adjacency failure occurs. In this case, $X \perp_1 Y$ and $Y \perp_1 Z$. Triple interactions that match the definition of 2-associations, however, need to have a specific structure. As we saw in Theorem 1, a 2-association always has to contain a collider.

From the above corollary, it is easy to see that in the chain graph $X \rightarrow Y \rightarrow Z$ we cannot find a 2-association, since none of the nodes is a collider. On the other hand,

the minimum unfaithful triple described in Example 1 matches the definition, since $X \overset{s}{\perp}_2 \{Y, Z\}$, $Y \overset{s}{\perp}_2 \{X, Z\}$ and $Z \overset{s}{\perp}_2 \{X, Y\}$. In general, strict 2-associations describe collider structures such as $X \rightarrow Y \leftarrow Z$ for which at least one of the edges violates adjacency faithfulness. We use this definition to introduce our new assumption.

Definition 8 (2-Adjacency Faithfulness) *Given $\mathcal{M} = (G, V, P)$, for all $X, Y \in V$, where X and Y are adjacent in the generating DAG G , there exists $Y \subseteq MB(X)$, with $Y \in Y$, s.t. $X \perp_{\leq 2} Y$.*

The main idea is to weaken adjacency faithfulness such that if a marginal dependence is not present, there will, however, be a dependence in combination with a parent, child or spouse. If we allow only 1-associations, 2-adjacency faithfulness would reduce to adjacency faithfulness. By also considering 2-associations, however, we can discover a larger spectrum of causal mechanisms.

The textbook example for a mechanism that violates faithfulness but is detectable by assuming 2-adjacency faithfulness is the xor-connection described in Example 1. Here, $Y \overset{s}{\perp}_2 \{X, Z\}$, two parents, while $X \overset{s}{\perp}_2 \{Y, Z\}$ —i.e. a child and a spouse. We could even slightly adapt the mechanism and only model Z using an unbiased coin but use a biased coin for X . In this case, only X is marginally independent of Y , while Z becomes dependent on Y . Moreover, assuming 2-adjacency faithfulness could even allow us to detect some faithfulness violations that are due to cancelling paths. In particular, consider the two paths $X \rightarrow Y$ and $X \rightarrow Z \rightarrow W \rightarrow Y$ that cancel such that $X \perp_P Y$. Since $X \perp_P Y$, $X \perp_P W \mid Z$, X could be strictly 2-associated to the set $\{W, Y\}$. Since we know that a 2-association contains a collider and we can neither find a 1-association to Y or W , we know that there has to be an edge violating adjacency faithfulness.

It is not possible to rely on orientation faithfulness when dealing with strict 2-associations. Although we know that a strict 2-association has to contain a collider, we do not know the skeleton structure within the triple and hence cannot apply orientation faithfulness. Next, we show that we can sometimes identify the collider if such a triple occurs in a larger graph.

6 ORIENTATION

So far, we showed how we can detect unfaithful triples from conditional (in)dependence statements under the weaker assumption of 2-adjacency faithfulness. Now imagine that we want to use this knowledge for causal discovery. If we observe an isolated triple that follows the dependence structure of the noisy xor, we can only tell that there is a collider. However, if we are given more

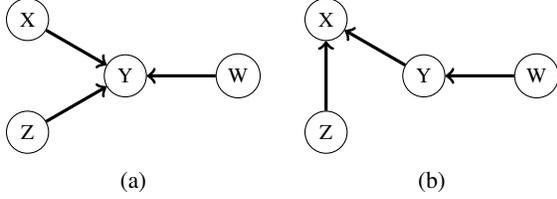


Figure 4: In both distributions $Y \perp_{-2} \{X, Z\}$ and $Y \perp_{-1} W$. In the graph shown in (a) Y is a collider on all paths between $\{X, Z\}$ and W , whereas in (b) Y is a non-collider.

information, we are able to break this symmetry.

Example 2 Consider that X and Z are unbiased coins as in the noisy xor example. In addition, there is a binary variable W with $P(W = 1) = p$, where $0 < p < 1$ and an unobserved binary noise variable E with $P(E = 1) = q$, where $0 < q < \frac{1}{2}$. Now we generate Y as

$$Y := ((X \oplus Z) \wedge W) \oplus E,$$

where E, W, X and Z are drawn independently. The requirements for p ensure that W is dependent on Y and the requirements on E ensure that the dependencies are non-deterministic ($q \neq 0$) and evident without observing E ($q \neq \frac{1}{2}$). The corresponding causal graph is given in Figure 4(a). From the induced dependencies, that we derive in detail in Supplementary Material S.1, we can now obtain an asymmetry. In particular, $\{X, Y, Z\}$ form a minimal unfaithful triple, but only Y is dependent on W , whereas $\{X, Z\} \perp_{-P} W$ and due to the xor, $X \perp_{-P} W \mid Y$ as well as $Z \perp_{-P} W \mid Y$. Thus, we can detect that there is no edge between X and W or Z and W since none of these pairs can be 2-associated. However, we do find that $X \not\perp_{-P} W \mid \{Y, Z\}$ and $Z \not\perp_{-P} W \mid \{Y, X\}$. As we will show in Theorem 2, we can use this information to identify that Y is the collider in the triple and that $W \rightarrow Y$.

To detect such an asymmetry, it is necessary that the collider in the triple is the effect of another node or pair of nodes. If, for example, X would be the collider in the triple and $W \rightarrow Y$ (see Figure 4(b)), we cannot find such an asymmetry. To generate that graph we could model Y as a noisy copy of W and construct X with a noisy xor from Y and Z . We still know that W is adjacent to Y , but we cannot direct any of the edges as for example we would find that $X \not\perp_{-P} W \mid Z$, which we would also observe if Z would be the collider in the triple, or if we would flip the edge direction between Y and W —i.e. W would be a noisy copy of Y .

Based on this intuition, we propose an orientation rule that may include causal structures that induce strict 2-associations. To do so, we use a shorthand notation—i.e. write $Y \rightarrow X$, if for each element $Y \in Y$ it holds that

$Y \rightarrow X$ and vice versa write $X \rightarrow Y$ if X is a parent of each node $Y \in Y$, that is, $\forall Y \in Y : X \rightarrow Y$.

Definition 9 (Orientation Rule) Let $M := (G, V, P)$ and we are given two disjoint sets $X, Z \subseteq V$ and $Y \in V$, where $Y \perp_{-2} X$ and $Y \perp_{-2} Z$, and no node $X \in X$ is adjacent to some node $Z \in Z$.

- i) If for each pair $X \in X$ and $Z \in Z$ it holds that X is dependent on Z given any subset of $V \setminus \{X, Z\}$ that contains $Y \cup (X \setminus \{X\}) \cup (Z \setminus \{Z\})$, then $X \rightarrow Y \leftarrow Z$,
- ii) otherwise, if for each pair $X \in X$ and $Z \in Z$ it holds that X is dependent on Z conditional on any subset of $V \setminus \{X, Z\}$ that contains $(X \setminus \{X\}) \cup (Z \setminus \{Z\})$ but does not contain Y , Y is a non-collider on at least one path $\langle X, Y, Z \rangle$ where $X \in X$ and $Z \in Z$.

Simply put, the above orientation rule relies on the fact that a 2-association contains a collider. Either Y is the collider on each path $\langle X, Y, Z \rangle$ between any variable $X \in X$ and $Z \in Z$ or Y is one of the parents in at least one of the triples and hence blocks at least one such path. If both sets X and Z only contain a single element, rule i) refers to a “normal” collider e.g. $X \rightarrow Y \leftarrow Z$ and rule ii) refers either to a chain like $X \rightarrow Y \rightarrow Z$ or to a common cause $X \leftarrow Y \rightarrow Z$. Let us consider Example 2 again, where we generated Y as a non-deterministic function of X, Z and W . First, we find that $Y \perp_{-2} \{X, Z\}$, $Y \perp_{-1} W$ and W is not adjacent to X or Z (since W is not 2-associated to X or Z), which is required to apply our rule. Further, we can apply rule i) since W is dependent on X given any set that includes $\{Y, Z\}$ and W is dependent on Z given any set that includes $\{Y, X\}$. Hence, we can infer that $\{X, Z\} \rightarrow Y \leftarrow W$.

In the following we will first show that our orientation rule is sound—i.e. if rule i) or ii) can be applied, then we are sure we found the corresponding graph structure—and then analyze the inverse, that is, what assumptions need to hold such that the given graph structure implies the suggested dependence model.

Theorem 2 Assuming that the causal Markov condition holds, the orientation rule in Definition 9 is sound.

We provide the proof for Theorem 2 in Supplementary Material S.3. We show both rules by contraposition, that is, to show the implication in rule ii) holds, we prove that if the true structure is $X \rightarrow Y \leftarrow Z$ (exactly the structure not implied by rule ii)), we can always find a pair $X \in X$ and $Z \in Z$ such that X becomes independent of Z if we condition on a set that includes $(X \setminus \{X\}) \cup (Z \setminus \{Z\})$, but does not contain Y . Rule i) can be proven accordingly.

The question that remains is: Does the inverse always hold? For example, if the true graph contains a non-collider structure such as $X \rightarrow Y \rightarrow Z$, will we always

find that $X \not\perp_P Z$? The short answer is no. Already when we only assume adjacency faithfulness, it can happen that $X \perp_P Z$ although the true graph is $X \rightarrow Y \rightarrow Z$ and it holds that $X \not\perp_P Y$ and $Y \not\perp_P Z$, which is called failure of transitivity. More generally, assuming that orientation faithfulness holds, such failures will not occur. In the following, we extend this assumption to our setting.

Definition 10 (2-Orientation Faithfulness) Let $M := (G, V, P)$ and we are given two disjoint sets $X, Z \subseteq V$ and $Y \in V$, where $Y \prec_{\leq 2} X$ and $Y \prec_{\leq 2} Z$, and no node $X \in X$ is adjacent to some node $Z \in Z$.

- i) If $X \rightarrow Y \leftarrow Z$ is in G , then for each pair $X \in X$ and $Z \in Z$, X is dependent on Z given any subset of $V \setminus \{X, Z\}$ that contains $Y \cup (X \setminus \{X\}) \cup (Z \setminus \{Z\})$,
- ii) otherwise, for each pair $X \in X$ and $Z \in Z$, X is dependent on Z conditional on any subset of $V \setminus \{X, Z\}$ that contains $(X \setminus \{X\}) \cup (Z \setminus \{Z\})$, but not Y .

Equivalently to 2-adjacency faithfulness, 2-orientation faithfulness reduces to orientation faithfulness, if both sets X and Z only contain a single element. For orientation faithfulness, it has been shown that all failures can be detected under the assumption that adjacency faithfulness holds (Zhang and Spirtes, 2008). Sadly, an equally strong statement cannot be made for 2-adjacency faithfulness and 2-orientation faithfulness, as we discuss in the following subsection.

6.1 Failures of 2-Orientation Faithfulness

Without any assumptions, we can detect triples for which $Y \prec_{\leq 2} X$ and $Y \prec_{\leq 2} Z$, and know by assuming CMC that all 2-associations contain a collider. If further, all paths $\langle X, Y, Z \rangle$ with $(X, Z) \in X \times Z$ are unshielded, we can detect if any of the conditions in 2-orientation faithfulness fails. In particular, due to the soundness of our orientation rule, we would detect that none of the conditions in the orientation rule is satisfied if condition i) or ii) in 2-orientation faithfulness fails, as we show in Corollary 1.

Yet, we cannot detect all failures of 2-orientation faithfulness. That is due to the fact that we might not always be able to detect whether all paths $\langle X, Y, Z \rangle$ are unshielded. If there is a direct edge between X and Z , we will always find that those nodes are either 1-associated or there exists a third node U such that they are 2-associated (if 2-adjacency faithfulness holds). However, if we find a strict 2-association between X and $\{Z, U\}$ there is no guarantee that the path is shielded. In particular, if U is the collider between X and Z , the triple is unshielded; but if Z is the collider between X and U , the triple is shielded (see Figure 5). In a causal discovery algorithm, we could try to iteratively infer the DAG structure within such triples until we cannot apply the

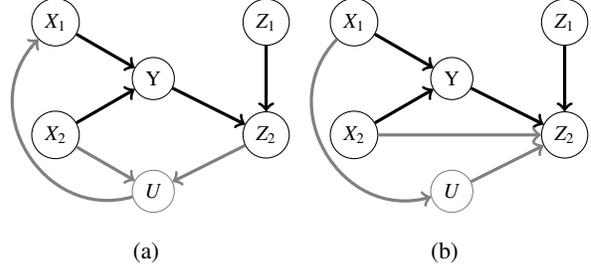


Figure 5: In both figures, $Y \prec_{\leq 2} X = \{X_1, X_2\}$, $Y \prec_{\leq 2} Z = \{Z_1, Z_2\}$ (related nodes and edges are marked in black) and $X_2 \prec_{\leq 2} \{U, Z_2\}$. If we are only given this information, we cannot determine whether the path $\langle X_2, Y, Z_2 \rangle$ is unshielded (a) or shielded (b). While in graph (a), we could safely apply our orientation rule, the shielded graph (b) can be problematic. Due to the directed path from X_1 over U to Z_2 and the adjacency between X_2 and Z_2 , each pair $X, Z \in X \times Z$ is now d -connected given $\{Y\} \cup (X \setminus \{X\}) \cup (Z \setminus \{Z\})$. Thus, the condition for rule i) could hold, although $X \rightarrow Y \leftarrow Z$ is not in G .

rule anymore. If we are lucky, we can first infer that $X \rightarrow U \leftarrow Z$ and after that also apply our rule for $\{X, Y, Z\}$. Keeping this exception in mind, we can derive the following corollary from Theorem 2.

Corollary 1 Given $M := (G, V, P)$ with $Y \in V$ and $X, Z \subseteq V$, where $X \cap Z = \emptyset$, $Y \prec_{\leq 2} X$, $Y \prec_{\leq 2} Z$ and no pair of nodes $(X, Z) \in X \times Z$ is adjacent. Assuming that CMC holds, we can detect if condition i) or ii) of 2-orientation faithfulness fails on the triple $\{X, Y, Z\}$.

The proof is provided in Supplementary Material S.3. In general, 2-orientation faithfulness might be useful not only for constraint-based causal discovery methods, but also for algorithms that aim to discover the Markov blanket of a target node or permutation-based causal discovery algorithms such as the Sparsest Permutation (SP) algorithm proposed by Raskutti and Uhler (2018). In Supplementary Material S.2, we provide a short discussion from which we conjecture that the SP algorithm can identify the collider pattern even for 2-associations like in Figure 4(a) if 2-orientation faithfulness holds.

In the next section, we demonstrate how to put theory into practice and propose an algorithm to find the Markov blanket of a target node under 2-adjacency faithfulness.

7 IMPLEMENTATION

As a proof of concept, we propose a simple modification of the Grow and Shrink (GS) algorithm (Margaritis and Thrun, 2000) to discover Markov blankets that can

contain 2-associations. After that, we briefly discuss further challenges that need to be solved to propose a causal discovery algorithm based on our new assumptions.

The GS algorithm is a simple and theoretically sound causal discovery algorithm, that as a first step identifies the Markov blanket for each node (Margaritis and Thrun, 2000). This step of the algorithm consists of a grow phase, in which we iteratively discover a superset of the Markov blanket of a target node Y , and a shrink phase, in which superfluous nodes are pruned.

To make sure that we can detect Markov blankets that contain strict 2-associations, we assume that 2-adjacency faithfulness holds and that we can detect all spouses. For the latter, we assume that a slight variation of the 2-orientation faithfulness assumption holds. In essence, we need to assume that condition i) in 2-orientation faithfulness also holds for shielded triples, which boils down to assuming that the spouses of the target do not cancel each other out, as we explain below.

Assumption 1 *Let $M := (G, V, P)$ and we are given two disjoint sets $X, Z \subseteq V$ and $Y \in V$, where $Y \prec_{\leq 2} X$ and $Y \prec_{\leq 2} Z$. If $X \rightarrow Y \leftarrow Z$ in G , then for each pair $X \in X$ and $Z \in Z$, X is dependent on Z given any subset of $V \setminus \{X, Z\}$ that contains $Y \cup (X \setminus \{X\}) \cup (Z \setminus \{Z\})$.*

The above assumption is a relatively lightweight adaption of condition i) in 2-orientation faithfulness. In particular, let $X = \{X, T\}$, where T is the target node. Then all nodes in Z are spouses of T and even become part of $PC(T)$ if all paths $\langle T, Y, Z \rangle$ for $Z \in Z$ are shielded. Thus, we would already add those nodes when looking for the parents and children of T . The only complication that may arise is if the second node $X \in X$ is adjacent to a node in $Z \in Z$ and this adjacency would lead to a cancellation such that Z is only dependent on T if we do not condition on X . The corresponding causal graph consists of the paths $T \rightarrow Y \leftarrow Z$ and $Y \leftarrow X \rightarrow Z$. Since X cannot block the path $\langle T, Y, Z \rangle$, such a scenario seems to be only possible if the causal mechanism that generates Z from X is deterministic. Based on this assumption, we can introduce our adapted GS algorithm.

The generalized GS algorithm is shown in Algorithm 1, where we only modified the grow phase to also consider pairs of random variables. This allows us to find nodes to which the target node is 2-associated or spouses to which a child node of T is 2-associated using Assumption 1. The shrink phase is not modified and checks if singletons can be removed. It is important to note that we will not remove single nodes of a true 2-association to T or a child of T , because we do not check for marginal dependencies. For example, assume that the pair $\{X, Z\}$ was added in the grow phase, where X is a child of T and Z

Algorithm 1: Modified GS for Markov Blankets

input : Random variables V with joint distribution P , Target $T \in V$

output: $MB(T)$

```

1  $V' \leftarrow V \setminus \{T\}$ ;
2  $S \leftarrow \emptyset$ ;
  // Grow Phase
3 while  $(\exists X \in V' : T \not\perp_P X \mid S) \vee$ 
4  $(\exists X, Z \in V' : T \not\perp_P X \mid S \cup \{Z\})$  do
5    $S \leftarrow S \cup \{X\}$ ;
  // Shrink Phase
6 while  $\exists X \in S : T \perp_P X \mid S \setminus X$  do
7    $S \leftarrow S \setminus X$ ;
8 return  $S$ 

```

the corresponding spouse. If we try to remove X in the shrink phase, we have that $T \not\perp_P X \mid S \setminus X$, since $Z \in S$. Hence, X remains in S , as well as Z .

In the following, we show that our proposed algorithm correctly identifies the Markov blanket of a target node assuming that 2-adjacency faithfulness, the causal Markov condition and Assumption 1 hold.

Theorem 3 *Given $M = (G, V, P)$. Assuming that 2-adjacency faithfulness, Assumption 1 and CMC hold, Algorithm 1 correctly identifies $MB(T)$ for $T \in V$.*

We provide the proof in Supplementary Material S.3. For discovering the Markov blanket, we do not need to know the collider of a 2-association since it only returns a set of nodes. The more challenging task is to implement our framework to discover causal networks. As an example, the next step in the GS algorithm would be to distinguish the spouses from the parents and children of a node. However, this is not straightforward for 2-associations, since we first need to identify the collider in the triple. Similarly, we could extend well-known algorithms such as the PC algorithm (Spirtes et al., 2000) or the GES algorithm (Chickering, 2002) by modifying the skeleton phase, respectively the forward phase such that we can find triple interactions as we did for GS. The edge orientation could be done by first applying the orientation rule in Definition 9 and then applying a similar set of rules like Meek’s orientation rules (Meek, 1995a). Alternatively, it was shown that SAT-based causal discovery algorithms can be easily adapted to weaker assumptions than faithfulness (Zhalama et al., 2017), which might be an interesting direction for future work.

8 CONCLUSION

In this work, we proposed 2-adjacency faithfulness, which is a weaker version of adjacency faithfulness. Our new assumption is able to detect faithfulness violations caused by weak or non-existent marginal dependencies, which are detectable by considering a combination of parents, children or spouses. We provide an in-depth analysis of such dependencies and propose a sound orientation rule, which can infer part of the correct causal structure by detecting colliders. We complement this rule with 2-orientation faithfulness, which assumes that if a causal graph contains such collider structures, we will find that the corresponding conditional dependence statements hold in P . As a proof of concept, we showed that we can extend the GS algorithm to find Markov blankets under strictly weaker assumptions than faithfulness.

For future work, we would like to develop a sound causal discovery algorithm based on 2-adjacency faithfulness and extend our theory to directed mixed graphs that can contain unobserved confounders.

Acknowledgements

This work was supported by the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement 639466). A. Marx is supported by the International Max Planck Research School for Computer Science (IMPRS-CS).

References

H. Andersen. When to Expect Violations of Causal Faithfulness and Why it Matters. *Philosophy of Science*, 80(5):672–683, 2013.

S. Bernstein. *Theory of Probability*. Moscow, 1927.

D. M. Chickering. Optimal Structure Identification With Greedy Search. *Journal of Machine Learning Research*, 3(Nov):507–554, 2002.

A. P. Dawid. Conditional Independence in Statistical Theory. *Journal of the Royal Statistical Society: Series B (Methodological)*, 41(1):1–15, 1979.

M. Forster, G. Raskutti, R. Stern, and N. Weinberger. The Frugal Inference of Causal Relations. *The British Journal for the Philosophy of Science*, 69(3):821–848, 2017.

D. Geiger, T. Verma, and J. Pearl. Identifying Independence in Bayesian Networks. *Networks*, 20(5):507–534, 1990.

G. Hesslow. Two Notes on the Probabilistic Approach to Causality. *Philosophy of science*, 43(2):290–292, 1976.

T. Inazumi, T. Washio, S. Shimizu, J. Suzuki, A. Yamamoto, and Y. Kawahara. Discovering causal structures in binary exclusive-or skew acyclic models. In *Proceedings of the International Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 373–382, 2011.

J. Lemeire, S. Meganck, F. Cartella, and T. Liu. Conservative independence-based causal structure learning in absence of adjacency faithfulness. *International Journal of Approximate Reasoning*, 53(9):1305–1325, 2012.

D. Margaritis and S. Thrun. Bayesian Network Induction via Local Neighborhoods. In *Advances in Neural Information Processing Systems*, pages 505–511, 2000.

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

C. Meek. Strong completeness and faithfulness in Bayesian networks. In *Proceedings of the International Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 411–419, 1995b.

J. Pearl. *Causality*. Cambridge university press, 2009.

J. Pearl and T. Verma. *The Logic of Representing Dependencies by Directed Graphs*. University of California (Los Angeles). Computer Science Department, 1987.

J. Peters, J. M. Mooij, D. Janzing, and B. Schölkopf. Causal discovery with continuous additive noise models. *Journal of Machine Learning Research*, 15(1):2009–2053, 2014.

J. Ramsey, P. Spirtes, and J. Zhang. Adjacency-Faithfulness and Conservative Causal Inference. In *Proceedings of the International Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 401–408. AUAI Press, 2006.

G. Raskutti and C. Uhler. Learning directed acyclic graph models based on sparsest permutations. *Stat*, 7(1), 2018.

D. Sejdinovic, A. Gretton, and W. Bergsma. A Kernel Test for Three-Variable Interactions. In *Advances in Neural Information Processing Systems*, pages 1124–1132, 2013.

P. Spirtes and J. Zhang. A Uniformly Consistent Estimator of Causal Effects under the k -Triangle-Faithfulness Assumption. *Statistical Science*, 29(4):662–678, 2014.

P. Spirtes, C. N. Glymour, R. Scheines, D. Heckerman, C. Meek, G. Cooper, and T. Richardson. *Causation, Prediction, and Search*. MIT press, 2000.

- W. Spohn. Stochastic Independence, Causal Independence, and Shieldability. *Journal of Philosophical Logic*, 9(1):73–99, 1980.
- C. Uhler, G. Raskutti, P. Bühlmann, and B. Yu. Geometry of the Faithfulness Assumption in Causal Inference. *The Annals of Statistics*, pages 436–463, 2013.
- T. Verma and J. Pearl. Equivalence and synthesis of causal models. In *Proceedings of the International Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 220–227, 1991.
- Zhalama, J. Zhang, F. Eberhardt, and W. Mayer. SAT-Based Causal Discovery under Weaker Assumptions. In *Proceedings of the 33th Annual Conference on Uncertainty in Artificial Intelligence (AUAI)*, 2017.
- J. Zhang. A Comparison of Three Occam’s Razors for Markovian Causal Models. *The British journal for the philosophy of science*, 64(2):423–448, 2013.
- J. Zhang and P. Spirtes. Detection of Unfaithfulness and Robust Causal Inference. *Minds and Machines*, 18(2): 239–271, 2008.
- J. Zhang and P. Spirtes. The three faces of faithfulness. *Synthese*, 193(4):1011–1027, Apr 2016.

SUPPLEMENTARY MATERIAL

S.1 Example 2 in Detail

As described in Section 6, we can generate a DAG of the form $X \rightarrow Y \leftarrow Z$ and $W \rightarrow Y$ s.t. X, Y and Z form a minimal unfaithful triple and $W \not\perp_P Y$ as follows. We generate X, Z, W and E independently, with X and Z as fair coins, W as a coin with $P(W = 1) = p$, where $0 < p < 1$ and E (the noise variable) as a biased coin with $P(E = 1) = q$, $0 < q < \frac{1}{2}$. With $q > 0$, we ensure that the function is non-deterministic. Further, we generate Y as

$$Y := ((X \oplus Z) \wedge W) \oplus E .$$

We will obtain that $P(Y = 1) = q + \frac{p}{2} - pq$. Further, we can calculate that $P(X = 1, Y = 1) = \frac{1}{2}P(Y = 1) = P(X = 1) \cdot P(Y = 1)$. Also, $P(X = 1, Y = 0) = P(X = 1) \cdot P(Y = 0)$, which means that they are marginally independent. The same holds for Z and Y . If we calculate the probability for all three variables, we get that $P(X = 0, Z = 1, Y = 1) = \frac{p+q-2pq}{4}$ and $P(X = 0, Z = 1) \cdot P(Y = 1) = \frac{1}{4}P(Y = 1)$. Hence, we need to solve

$$P(X = 0, Z = 1, Y = 1) = P(X = 0, Z = 1) \cdot P(Y = 1)$$

$$\Leftrightarrow p + q - 2pq = q + \frac{p}{2} - pq$$

$$\Leftrightarrow p - pq = \frac{p}{2} .$$

The only solutions are $p = 0$ or $q = \frac{1}{2}$, which we excluded. Hence, $Y \not\perp_P \{X, Z\}$ and by weak union also $Y \not\perp_P X | Z$, as well as $Y \not\perp_P Z | X$. Since we know by assumption that $X \perp_P T$ we can conclude from Lemma 1 that also $X \not\perp_P Z | Y$, which means that $\{X, Y, Z\}$ from a minimal unfaithful triple since W will also not cancel out any of these conditional dependencies. Next, we also find that $W \not\perp_P Y$, since $P(W = 1, Y = 1) = \frac{p}{2}$, which is only equal to $P(W = 1) \cdot P(Y = 1)$, if $p = 0$, $p = 1$ or $q = \frac{1}{2}$, which we excluded, and hence $W \not\perp_P Y$. Last, we need to show that $X \not\perp_P W | \{Y, Z\}$ and that $Z \not\perp_P W | \{X, Y\}$. We can write

$$P(X, W | Y, Z) = \frac{P(X, W, Y, Z)}{P(Y, Z)} .$$

To show conditional dependence, this value has to be different from $P(X | Y, Z) \cdot P(W | Y, Z)$. Consider the case where all variables are equal to one. Hence, we get that

$$P(X = 1, W = 1, Y = 1, Z = 1) = \frac{pq}{4} ,$$

$$P(X = 1, Y = 1, Z = 1) = \frac{q}{4} ,$$

$$P(W = 1, Y = 1, Z = 1) = \frac{p}{4} .$$

Since we know that $P(Y = 1, Z = 1) = P(Y = 1)/2$, we thus need to solve

$$pq = \frac{pq}{2P(Y = 1)} .$$

This equation can only be true if p or $q = 0$, i.e. the system is either independent of W or deterministic, $p = 1$ or $q = \frac{1}{2}$, which we all excluded by assumption. Hence, $X \not\perp_P W | \{Y, Z\}$. The dependence between Z and W given X and Y can be derived in the same way.

S.2 2-Orientation Faithfulness and Sparsest Markov Representation

In this section, we briefly discuss the connection of our new assumptions to approaches based on the sparsest Markov representation (SMR) (Raskutti and Uhler, 2018) which is also referred to as frugality (Forster et al., 2017), which we discussed in the related work section. A graph G^* satisfies the SMR assumption if every graph G that fulfils the Markov property and is not in the Markov equivalence class of G^* contains more edges than G^* . Here we will not discuss the SMR assumption in further detail, but focus on the suggested causal discovery algorithm under the SMR assumption, which is called the Sparsest Permutation (SP) algorithm.

To explain the SP algorithm, we need to define a DAG G_π , w.r.t. a permutation π . A DAG G_π consists of vertices V and directed edges E_π , where an edge from the j -th

node $\pi(j)$ according to permutation π to node $\pi(k)$ is in E_π if and only if $j < k$ and

$$X_{\pi(j)} \not\perp_P X_{\pi(k)} \mid \{X_{\pi(1)}, X_{\pi(2)}, \dots, X_{\pi(k-1)}\} \setminus \{X_{\pi(j)}\},$$

where $X_{\pi(j)}$ refers to the j -th random variable according to permutation π . Based on this definition, the SP algorithm constructs a graph G_π for each possible permutation and selects that permutation π^* for which G_{π^*} contains the fewest edges. This permutation π^* is also called minimal or a minimal permutation, if it is not unique.

Although this procedure might be very slow in practice, it has theoretically appealing properties. In particular, we conjecture that it can identify the collider pattern even if strict 2-associations are included, if 2-orientation faithfulness holds. In this work, we will not provide a proof for this conjecture, but give some evidence by discussing the behaviour of the SP algorithm on an example graph.

Consider the graph provided in Figure 4(a) again. For this example, we assume that V does not consist of any further vertices than the four shown in the graph. We will show that all permutations π that are minimal have in common that $\pi(4) = Y$. W.l.o.g. let $\pi(1) = X, \pi(2) = Z$ and $\pi(3) = W$, then G_π only contains the four correct edges, which are:

$$\begin{aligned} \pi(1) \rightarrow \pi(4) &: X \not\perp_P Y \mid \{Z, W\} \\ \pi(2) \rightarrow \pi(4) &: Z \not\perp_P Y \mid \{X, W\} \\ \pi(3) \rightarrow \pi(4) &: W \not\perp_P Y \mid \{X, Z\} \end{aligned}$$

and we do not add any superfluous edges, as

$$\begin{aligned} \pi(1) \rightarrow \pi(2) &: X \perp_P Z \mid \emptyset \\ \pi(1) \rightarrow \pi(3) &: X \perp_P W \mid Z \\ \pi(2) \rightarrow \pi(3) &: Z \perp_P W \mid X. \end{aligned}$$

If we would pick a permutation π' in which we flip for example W and Y such that Y is no longer the node assigned to the highest number in the permutation, i.e. $\pi'(3) = Y$ and $\pi'(4) = W$, we will find more edges and thus not a minimal graph anymore. In particular, we get that

$$\begin{aligned} \pi'(1) \rightarrow \pi'(3) &: X \not\perp_P Y \mid \{Z\} \\ \pi'(2) \rightarrow \pi'(3) &: Z \not\perp_P Y \mid \{X\} \\ \pi'(3) \rightarrow \pi'(4) &: Y \not\perp_P W \mid \{X, Z\} \\ \pi'(1) \rightarrow \pi'(4) &: X \not\perp_P W \mid \{Z, Y\} \\ \pi'(2) \rightarrow \pi'(4) &: Z \not\perp_P W \mid \{X, Y\} \end{aligned}$$

and thus the graph according to this permutation contains two edges more than for permutation π . The main

point is that we are now allowed to condition on Y , which opens the paths between X or Z and W . Similarly, assume that we put X as the last node and get the order $\pi'(1) = Z, \pi'(2) = W, \pi'(3) = Y$ and $\pi'(4) = X$, for which

$$\begin{aligned} \pi'(1) \rightarrow \pi'(2) &: Z \perp_P W \mid \emptyset \\ \pi'(1) \rightarrow \pi'(3) &: Z \perp_P Y \mid \{W\} \\ \pi'(1) \rightarrow \pi'(4) &: Z \not\perp_P X \mid \{W, Y\} \\ \pi'(2) \rightarrow \pi'(3) &: W \not\perp_P Y \mid \{Z\} \\ \pi'(2) \rightarrow \pi'(4) &: W \not\perp_P X \mid \{Z, Y\} \\ \pi'(3) \rightarrow \pi'(4) &: Y \not\perp_P X \mid \{Z, W\} \end{aligned}$$

and hence, we again find four edges, which is one more than for π . Also, if $\pi'(1) = Y$, we can use it in the conditional to find a dependence between X and Z and at least one dependence between X or Z and W . Hence, at least for this example graph, the SP algorithm would infer a correct ordering.

An interesting avenue for future work would be to analyze whether it is possible to always detect the collider pattern also in larger graphs and triples that may or may not be shielded.

S.3 Additional Proofs

Theorem 2 *Assuming that the causal Markov condition holds, the orientation rule in Definition 9 is sound.*

Proof: First, we derive a general statement about the relations between X and Z without further specifying the role of Y . In particular, we show that there always exists a pair $(X, Z) \in X \times Z$ s.t. w.l.o.g.

$$X \perp_G Z \mid \text{Pa}(X) \cup (X \setminus \{X\}) \cup (Z \setminus \{Z\}), \quad (1)$$

where $\text{Pa}(X) \subseteq V \setminus Z$. Due to acyclicity, there has to exist a node in $X \cup Z$, say X , that is not an ancestor of any node in $(X \cup Z) \setminus \{X\}$ and hence $(X \cup Z) \setminus \{X\} \subseteq \text{Nd}(X)$. By the local Markov condition, we get that $X \perp_G (X \cup Z) \setminus \{X\} \mid \text{Pa}(X)$. Thus, by weak union,

$$X \perp_G Z \mid \text{Pa}(X) \cup (X \setminus \{X\}) \cup (Z \setminus \{Z\}),$$

for any $Z \in Z$. Further, $Z \cap \text{Pa}(X) = \emptyset$, as by assumption no pair of nodes $(X, Z) \in X \times Z$ is adjacent in G .

Since $Y \prec_{\leq} X$ and $Y \prec_{\leq} Z$, we know that Y is at least adjacent to one node in X and one node in Z . Hence, Y can take the following roles:

- Y is a descendent of each node in $X \cup Z$ (which corresponds to $X \rightarrow Y \leftarrow Z$),
- Y is a non-descendent of each node in $X \cup Z$ and

- c) Y is a descendent of at least one node in $X \cup Z$ and a non-descendent of at least one node in $X \cup Z$.

The first statement corresponds to the graph structure implied by rule i) and any possible structure from the latter two is implied by the probabilities found in rule ii). To show these two implications hold, we do a proof by contraposition for each rule.

Hence, to show rule i), we need to prove that if the graph structure is not a collider—i.e. Y takes one of the roles described in b) or c)—then there exists a pair $(X, Z) \in X \times Z$ and there exists a subset $S \subseteq V \setminus \{X, Z\}$ s.t.

$$X \perp\!\!\!\perp_P Z \mid S \cup \{Y\} \cup (X \setminus \{X\}) \cup (Z \setminus \{Z\}).$$

First, consider all graphs in which Y is a non-descendent of each node in $X \cup Z$ as described in b) We know from statement (1) that, w.l.o.g., there exists a pair $(X, Z) \in X \times Z$ for which $X \perp\!\!\!\perp_G Z \mid \text{Pa}(X) \cup (X \setminus \{X\}) \cup (Z \setminus \{Z\})$. Since $Y \in \text{Nd}(X)$, we will also find that $X \perp\!\!\!\perp_G Z \mid \text{Pa}(X) \cup (X \setminus \{X\}) \cup (Z \setminus \{Z\}) \cup \{Y\}$, where $\text{Pa}(X)$ does not include X or Z . Thus, by CMC we found the required independence. For the cases described in c), again assume that X is not an ancestor of any node in $(X \cup Z) \setminus \{X\}$. To conclude the same statement as previously, we show that X has to be in $\text{De}(Y)$ and thus $Y \in \text{Nd}(X)$. We do this by deriving a contradiction: assume $X \in \text{Nd}(Y)$. If X consists only of the single node X , then X has to be adjacent to Y , $X \in \text{Pa}(Y)$ and hence $X \rightarrow Y$ in G . Thus, Y (and hence X) has to be an ancestor of at least one node in Z , by assumption (Y is a non-descendent of at least one node in $X \cup Z$), which is a contradiction. Similarly, if X contains a second node, X' , we know by assumption that $X' \in \text{Nd}(X)$. We also know that the triple $\{X, X', Y\}$ has to contain a collider. X cannot be the collider, since $X \notin \text{De}(Y)$ and also X' cannot be the collider since $X \notin \text{An}(X')$. Hence, Y has to be the collider on the path $\langle X, Y, X' \rangle$. As above, at least one node $Z \in Z$ has to be a descendent of Y , by assumption and thus, $X \in \text{An}(Z)$, which is a contradiction.

Last, we prove that the implication in rule ii) holds. Thus, by contraposition, we need to show that if $X \rightarrow Y \leftarrow Z$, then there exists a pair $X, Z \in X \times Z$ s.t. X is conditionally independent of Z given a subset of $V \setminus \{X, Z\}$ that contains $(X \setminus \{X\}) \cup (Z \setminus \{Z\})$ but does not contain Y . From statement (1) there exists a pair $(X, Z) \in X \times Z$ that is d -separated given $\text{Pa}(X) \cup (X \setminus \{X\}) \cup (Z \setminus \{Z\})$. Since Y cannot be in $\text{Pa}(X)$ due to acyclicity, we showed that there exists such a pair of nodes X, Z that can be rendered conditionally independent by a subset of $V \setminus \{X, Z\}$ that contains $(X \setminus \{X\}) \cup (Z \setminus \{Z\})$ but does not contain Y (after applying CMC), which concludes the proof. \square

Corollary 1 Given $M := (G, V, P)$ with $Y \in V$ and $X, Z \subseteq V$, where $X \cap Z = \emptyset$, $Y \prec_{\leq} X$, $Y \prec_{\leq} Z$ and

no pair of nodes $(X, Z) \in X \times Z$ is adjacent. Assuming that CMC holds, we can detect if condition i) or ii) of 2-orientation faithfulness fails on the triple $\{X, Y, Z\}$.

Proof: Since we know that $Y \prec_{\leq} X$ and $Y \prec_{\leq} Z$, we can conclude that, as in the proof of Theorem 2, Y can take three different roles w.r.t. X and Z , where role a) corresponds to condition i) in 2-orientation faithfulness and rule i) in the orientation rule and roles b) and c) correspond to condition ii) and rule ii).

Now assume that condition i) in 2-orientation faithfulness fails, that is, the true graph can be described by role a), but there exists a pair $X \in X$ and $Z \in Z$, for which X is independent of Z given a subset of $V \setminus \{X, Z\}$ that contains $Y \cup (X \setminus \{X\}) \cup (Z \setminus \{Z\})$. If this is the case, we cannot apply rule i) of our orientation rule. In addition, we showed in Theorem 2 that for a graph as described by a) rule ii) can never apply. Thus, we can detect this failure of condition i) in 2-orientation faithfulness by noticing that neither rule i) nor ii) of our orientation rule applies.

Next, assume condition ii) in 2-orientation fails. This means that we cannot apply rule ii) of the orientation rule. Again, we showed that for such graphs Y takes either role b) or c), in which case orientation rule i) can never apply. Hence, we can detect if condition ii) in 2-orientation faithfulness fails, since none of the two conditions in our orientation rule are met. \square

Theorem 3 Given $M = (G, V, P)$. Assuming that 2-adjacency faithfulness, Assumption 1 and CMC hold, Algorithm 1 correctly identifies $\text{MB}(T)$ for $T \in V$.

Proof: We follow the original correctness proof under the faithfulness assumption (Margaritis and Thrun, 2000), that consists of two main steps. First, we need to show that $\text{MB}(T) \subseteq S$ after the grow phase and second, we need to ensure that all nodes in $\text{MB}(T)$ stay in S during the shrink phase, while nodes not in $\text{MB}(T)$ will be removed from S in the shrink phase.

Grow phase: By assumption (2-adjacency faithfulness), for each node $X \in \text{PC}(T)$, T is either 1-associated to X , or there exists a set X that includes X such that $T \rightarrow_2 X$. If T is 1-associated to a node X , then $T \not\perp\!\!\!\perp_P X \mid S$, if $X \notin S$, hence we will add those nodes. If T is 2-associated to a set $\{X, Z\}$ then $T \not\perp\!\!\!\perp_P X \mid S \cup \{Z\}$ for all $S \subseteq V \setminus \{X, T, Z\}$. Thus, we also add X to S , if $X \notin S$ and afterwards also find that $T \not\perp\!\!\!\perp_P Z \mid S$, if $Z \notin S$, since $X \in S$. Hence, all nodes in $\text{PC}(T)$ will be added during the grow phase. Next, we need to consider the spouses of T that do not overlap with $\text{PC}(T)$, hence might not have been added yet.⁵ Since we know that eventually S will contain all

⁵There could be nodes that are spouses of T and in $\text{PC}(T)$

children of T , we will afterwards also add the corresponding spouses. In particular, we need to consider two classes of spouses S : 1) Spouses that through a child node C are 2-associated to T ($T \dashv_2 \{C, S\}$). Those will be added due to the 2-association as explained above. 2) Spouses that are not involved in such a 2-association. For the latter, we find a conditional dependence between T and S by conditioning on the corresponding child node C (by Assumption 1), which will be in \mathcal{S} . A special case occurs if a child node C is 2-associated to two spouses S_1 and S_2 . Due to Assumption 1, T is dependent on S_1 if we condition on C and S_2 , vice versa T is dependent on S_2 if we condition on C and S_1 . Similarly to how we add 2-associations above, we will also first add one of the two and then the second one. Thus, after the grow phase, \mathcal{S} will contain all elements of $\text{MB}(T)$.

Shrink phase: Since it is possible that after the grow phase \mathcal{S} is a superset of $\text{MB}(T)$, we need to ensure that in the shrink phase all $W \notin \text{MB}(T)$ will be deleted from \mathcal{S} and all $X \in \text{MB}(T)$ will stay in \mathcal{S} .

First, we show that no node $X \in \text{MB}(T)$ will be removed from \mathcal{S} . Assume X is the first element in $\text{MB}(T)$ that we attempt to remove from \mathcal{S} . If $X \in \text{PC}(T)$, by definition of 2-adjacency faithfulness T is either 1-associated to X and hence, X will not be removed, or T is 2-associated with a set $X \subseteq \text{MB}(T)$ that contains X . W.l.o.g. let $X = \{X, Z\}$, then $T \not\perp_P X \mid \mathcal{S} \setminus \{X\}$, since \mathcal{S} contains Z , and hence, X will not be removed from \mathcal{S} . If X is a spouse of T , there again exist two cases. Either T is 2-associated to a set that contains X , in which case, X will not be removed from \mathcal{S} as explained above, or T is not 2-associated to a set that contains X . In the latter case, by Assumption 1, X is dependent on T conditioned on a subset of $\text{MB}(T) \setminus \{X\}$ and thus $X \not\perp_P T \mid \mathcal{S} \setminus \{X\}$. In particular, this subset consists of the common child C and in the special case that C is 2-associated to X and a second spouse S , it also contains that second spouse S . Either way, those conditioning sets are contained in \mathcal{S} . Hence, X will not be removed from \mathcal{S} . In the following iterations, \mathcal{S} will still contain $\text{MB}(T)$ and hence, we will also not remove a true element of $\text{MB}(T)$.

Last, assume $W \notin \text{MB}(T)$, but $W \in \mathcal{S}$ after the grow phase. Further, we can write $\mathcal{S} \setminus \{W\}$ as $\text{MB}(T) \cup \mathcal{Q}$, where \mathcal{Q} contains all elements from $\mathcal{S} \setminus \{W\}$ that are not in $\text{MB}(T)$. Then, $T \perp_G \{W\} \cup \mathcal{Q} \mid \text{MB}(T)$ and thus by weak union, $T \perp_G W \mid \text{MB}(T) \cup \mathcal{Q}$, which implies $T \perp_P W \mid \mathcal{S} \setminus \{W\}$ (by CMC). Hence, we delete each node in \mathcal{S} that is not in $\text{MB}(T)$ in the shrink phase. \square

at the same time e.g. if T has two children X and Z , where Z is also a parent of X .