

---

# Greedy Relaxations of the Sparsest Permutation Algorithm (Supplementary material)

---

Wai-Yin Lam<sup>1</sup>

Bryan Andrews<sup>1</sup>

Joseph Ramsey<sup>1</sup>

<sup>1</sup> Department of Philosophy, Carnegie Mellon University, Pittsburgh, Pennsylvania, USA

## A BACKGROUND MATERIALS

### A.1 GRAPHICAL DEFINITIONS

A *directed graph*  $\mathcal{G}$  over a set of measured variables  $\mathbf{V} = \{X_1, \dots, X_m\}$  consists of  $m$  vertices  $\mathbf{v} = \{1, \dots, m\}$  where each vertex  $i \in \mathbf{v}$  associates to a variable  $X_i \in \mathbf{V}$ , and each edge in  $\mathcal{G}$  is directed with the form  $j \rightarrow k$  and no vertex has a directed edge to itself. A *path*  $\mathbf{p}$  is a sequence of vertices  $\langle \mathbf{i}_1, \mathbf{i}_2, \dots, \mathbf{i}_k \rangle$  for some  $k \geq 2$  where  $\mathbf{i}_j \in \mathbf{v}$  for each  $1 \leq j \leq k$ , and  $\mathbf{i}_j$  and  $\mathbf{i}_{j+1}$  are connected by a directed edge (i.e.,  $\mathbf{i}_j \rightarrow \mathbf{i}_{j+1}$  or  $\mathbf{i}_{j+1} \rightarrow \mathbf{i}_j$ ). A path  $\mathbf{p}$  is *directed* if  $\mathbf{i}_j \rightarrow \mathbf{i}_{j+1}$  for each  $1 \leq j < k$ . A *directed acyclic graph* (DAG) is a directed graph where no vertex can have a directed path to itself.

Denote  $E(\mathcal{G})$  as the set of directed edges in  $\mathcal{G}$ . A pair of DAGs  $\mathcal{G}_1, \mathcal{G}_2$  over the same set of variables  $\mathbf{V}$  are equivalent if and only if  $E(\mathcal{G}_1) = E(\mathcal{G}_2)$ . Let  $\text{Pa}(j, \mathcal{G}) = \{k \in \mathbf{v} : (k \rightarrow j) \in E(\mathcal{G})\}$  be the set of *parents* of  $j$  in  $\mathcal{G}$ , and  $\text{Ch}(j, \mathcal{G}) = \{k \in \mathbf{v} : (j \rightarrow k) \in E(\mathcal{G})\}$  be the set of *children* of  $j$  in  $\mathcal{G}$ .  $\text{An}(j, \mathcal{G})$ , the *ancestors* of  $j$  in  $\mathcal{G}$ , is defined by the transitive closure of  $\text{Pa}(j, \mathcal{G})$ . Similarly,  $\text{De}(j, \mathcal{G})$ , the *descendants* of  $j$  in  $\mathcal{G}$ , is defined by the transitive closure of  $\text{Ch}(j, \mathcal{G})$  and union with  $\{j\}$  itself (i.e.,  $j$  is its own descendant). Further let  $\text{Nd}(j, \mathcal{G}) = \mathbf{v} \setminus \text{De}(j, \mathcal{G})$  be the set of  $j$ 's *non-descendants*.

A pair of vertices  $j, k \in \mathbf{v}$  are said to be *adjacent* in  $\mathcal{G}$  if  $(j \rightarrow k) \in E(\mathcal{G})$  or  $(k \rightarrow j) \in E(\mathcal{G})$ . For any triple of pairwise distinct vertices  $i, j, k \in \mathbf{v}$ , we say that  $(i, j, k)$  is *unshielded* if  $(i, j)$  and  $(j, k)$  are adjacent pairs in  $\mathcal{G}$ , but not  $(i, k)$ .  $(i, j, k)$  forms a *triangle* if they are pairwise adjacent. If  $(i, j, k)$  is an unshielded triple or is a triangle,  $j$  is a *collider* (on the path  $\langle i, j, k \rangle$ ) if  $(i \rightarrow j), (k \rightarrow j) \in E(\mathcal{G})$ , and a *non-collider* otherwise. A path  $\mathbf{p}$  is a *trek* if it contains no collider.

For any  $j, k \in \mathbf{v}$  and any  $\mathbf{i} \subseteq \mathbf{v} \setminus \{j, k\}$ ,  $j$  and  $k$  are *d-connected* given  $\mathbf{i}$  in  $\mathcal{G}$  if there exists a path  $\mathbf{p}$  between  $j$  and  $k$  in  $\mathcal{G}$  such that no non-collider on  $\mathbf{p}$  is in  $\mathbf{i}$ , and each collider  $l$  on  $\mathbf{p}$  or a  $l$ 's descendant is in  $\mathbf{i}$ .  $j$  and  $k$  are *d-separated* given  $\mathbf{i}$  in  $\mathcal{G}$  if  $j$  and  $k$  are not d-connected given  $\mathbf{i}$ . For any disjoint subsets of vertices  $\mathbf{j}, \mathbf{k}, \mathbf{i} \subseteq \mathbf{v}$ ,  $\mathbf{j}$  and  $\mathbf{k}$  are d-separated given  $\mathbf{i}$  in  $\mathcal{G}$  if  $j$  and  $k$  are d-separated by  $\mathbf{i}$  in  $\mathcal{G}$  for every  $j \in \mathbf{j}$  and every  $k \in \mathbf{k}$ .

Given a model  $(\mathcal{G}, \mathcal{P})$  over  $\mathbf{V}$ ,  $\mathcal{G}$  is said to be *local Markov* to  $\mathcal{P}$  if  $X_j \perp\!\!\!\perp_{\mathcal{P}} \mathbf{X}_{\text{Nd}(j, \mathcal{G})} \setminus \mathbf{X}_{\text{Pa}(j, \mathcal{G})} \mid \mathbf{X}_{\text{Pa}(j, \mathcal{G})}$  for every  $j \in \mathbf{v}$ . It is a well-known fact that  $\mathcal{G}$  is local Markov to  $\mathcal{P}$  if and only if  $I(\mathcal{G}) \subseteq I(\mathcal{P})$  (i.e., global Markov as defined by d-separation).

### A.2 GRAPHOID AXIOMS

For any pairwise disjoint sets of variables  $\mathbf{W}, \mathbf{X}, \mathbf{Y}$ , and  $\mathbf{Z}$ ,

$$\begin{aligned}
 \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z} &\Rightarrow \mathbf{Y} \perp\!\!\!\perp \mathbf{X} \mid \mathbf{Z} && \text{(symmetry)} \\
 \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \cup \mathbf{W} \mid \mathbf{Z} &\Rightarrow (\mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z}) \wedge (\mathbf{X} \perp\!\!\!\perp \mathbf{W} \mid \mathbf{Z}) && \text{(decomposition)} \\
 \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \cup \mathbf{W} \mid \mathbf{Z} &\Rightarrow \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z} \cup \mathbf{W} && \text{(weak union)} \\
 (\mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z}) \wedge (\mathbf{X} \perp\!\!\!\perp \mathbf{W} \mid \mathbf{Z} \cup \mathbf{Y}) &\Rightarrow \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \cup \mathbf{W} \mid \mathbf{Z} && \text{(contraction)} \\
 (\mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z} \cup \mathbf{W}) \wedge (\mathbf{X} \perp\!\!\!\perp \mathbf{W} \mid \mathbf{Z} \cup \mathbf{Y}) &\Rightarrow \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \cup \mathbf{W} \mid \mathbf{Z} && \text{(intersection)} \\
 (\mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{Z}) \wedge (\mathbf{X} \perp\!\!\!\perp \mathbf{W} \mid \mathbf{Z}) &\Rightarrow \mathbf{X} \perp\!\!\!\perp \mathbf{Y} \cup \mathbf{W} \mid \mathbf{Z} && \text{(composition)}
 \end{aligned}$$

A distribution  $\mathcal{P}$  is a *semigraphoid* if  $\mathbf{I}(\mathcal{P})$  is closed under *symmetry*, *decomposition*, *weak union*, and *contraction*. A semigraphoid  $\mathcal{P}$  is a *graphoid* if  $\mathbf{I}(\mathcal{P})$  is closed under *intersection*. A graphoid is *compositional* if  $\mathbf{I}(\mathcal{P})$  is closed under *composition*. See Chapter 2 of [Studený, 2005] for a more comprehensive study of graphoid axioms. In addition, applications of *symmetry* in our upcoming proofs will be done implicitly for the sake of simplicity.

Additionally, Spohn [1994] notes that the following property necessarily holds in the independence models induced by positive discrete probability distributions. For any pairwise disjoint sets of variables  $\mathbf{W}$ ,  $\mathbf{X}$ ,  $\mathbf{Y}$ , and  $\mathbf{Z}$ ,

$$(\mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid \mathbf{W} \cup \mathbf{Z}) \wedge (\mathbf{W} \perp\!\!\!\perp \mathbf{Z} \mid \mathbf{X} \cup \mathbf{Y}) \wedge (\mathbf{W} \perp\!\!\!\perp \mathbf{Z} \mid \mathbf{X}) \Rightarrow [(\mathbf{W} \perp\!\!\!\perp \mathbf{Z} \mid \mathbf{Y}) \Leftrightarrow (\mathbf{W} \perp\!\!\!\perp \mathbf{Z} \mid \emptyset)] \quad (\text{Spohn condition})$$

### A.3 DAG INDUCED FROM A PERMUTATION

**Definition A.1** Given a semigraphoid  $\mathcal{P}$  over  $\mathbf{V}$ , for every  $X \in \mathbf{V}$ , we say that  $\mathbf{M} \subseteq \mathbf{V}$  is a *Markov blanket* of  $X$  relative to  $\mathbf{Z} \subseteq \mathbf{V} \setminus \{X\}$  if

- (i)  $\mathbf{M} \subseteq \mathbf{Z}$ ;
- (ii)  $X \perp\!\!\!\perp_{\mathcal{P}} (\mathbf{Z} \setminus \mathbf{M}) \mid \mathbf{M}$ .

Such a Markov blanket  $\mathbf{M}$  is said to be a *Markov boundary* if it further satisfies the following condition:

- (iii) there does not exist  $\mathbf{M}' \subset \mathbf{M}$  s.t.  $X \perp\!\!\!\perp_{\mathcal{P}} (\mathbf{Z} \setminus \mathbf{M}') \mid \mathbf{M}'$ .

**Lemma A.2** [Verma and Pearl, 1988] Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , for every  $X \in \mathbf{V}$  and every  $\mathbf{Z} \subseteq \mathbf{V} \setminus \{X\}$ , there is a *unique Markov boundary* of  $X$  relative to  $\mathbf{Z}$ .

In the following, we use  $\text{MB}_{\mathcal{P}}(X, \mathbf{Z})$  to refer to the unique Markov boundary of  $X$  relative to  $\mathbf{Z}$ . The subscript  $\mathcal{P}$  will be suppressed if the underlying graphoid is clear from context.

**Lemma A.3** Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , for every  $X \in \mathbf{V}$  and every  $\mathbf{Z} \subseteq \mathbf{V} \setminus \{X\}$ , if  $\mathbf{M}$  is a Markov blanket of  $X$  relative to  $\mathbf{Z}$ , then  $\text{MB}(X, \mathbf{Z}) \subseteq \mathbf{M}$ .

*Proof.* Immediate from **Definition A.1** and **Lemma A.2**. □

Next, we revisit the two methods of inducing a DAG from a permutation. Given a semigraphoid  $\mathcal{P}$  over  $\mathbf{V}$ , each  $\pi \in \Pi(\mathbf{v})$  induces a DAG satisfying the following condition:

$$X_j \in \mathbf{M} \Leftrightarrow (j \rightarrow k) \in \mathbf{E}(\mathcal{G}_{\pi}) \quad (\text{VP})$$

where  $\mathbf{M}$  is a Markov boundary of  $X_k$  relative to  $\mathbf{X}_{\text{Pre}(k, \pi)}$ . (VP) is the construction of a *boundary DAG* in [Verma and Pearl, 1988]. On the other hand, given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , each  $\pi \in \Pi(\mathbf{v})$  induces a DAG satisfying the following condition:

$$j \in \text{Pre}(k, \pi) \text{ and } X_j \not\perp\!\!\!\perp_{\mathcal{P}} X_k \mid \mathbf{X}_{\text{Pre}(k, \pi) \setminus \{j\}} \Leftrightarrow (j \rightarrow k) \in \mathbf{E}(\mathcal{G}_{\pi}). \quad (\text{RU})$$

We want to show that the two DAG-inducing methods are equivalent when the underlying distribution is a graphoid.

**Lemma A.4** Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , consider any  $\pi \in \Pi(\mathbf{v})$ . Let  $\mathcal{G}_{\pi}$  be the DAG induced from  $\pi$  by (VP), and  $\mathcal{H}_{\pi}$  be the DAG induced from  $\pi$  by (RU). Then  $\mathcal{G}_{\pi} = \mathcal{H}_{\pi}$ .

*Proof.* We divide the proof into two directions: (VP)  $\Rightarrow$  (RU), and (VP)  $\Leftarrow$  (RU). Consider any  $j, k \in \mathbf{v}$  where  $\pi[j] < \pi[k]$  such that  $j \in \text{Pre}(k, \pi)$ . Let  $\mathbf{M}$  be the unique Markov boundary  $\text{MB}(X_k, \mathbf{X}_{\text{Pre}(k, \pi)})$ .

[ $\Rightarrow$ ] Suppose that  $(j \rightarrow k) \notin \mathbf{E}(\mathcal{G}_{\pi})$ . We have  $X_j \notin \mathbf{M}$ . By **Definition A.1** (ii), we then have,

$$X_k \perp\!\!\!\perp_{\mathcal{P}} ((\mathbf{X}_{\text{Pre}(k, \pi)} \setminus \mathbf{M}) \setminus \{X_j\}) \cup \{X_j\} \mid \mathbf{M} \quad \because X_k \perp\!\!\!\perp_{\mathcal{P}} \mathbf{X}_{\text{Pre}(k, \pi)} \setminus \mathbf{M} \mid \mathbf{M} \quad (1)$$

$$X_k \perp\!\!\!\perp_{\mathcal{P}} X_j \mid \mathbf{M} \cup ((\mathbf{X}_{\text{Pre}(k, \pi)} \setminus \mathbf{M}) \setminus \{X_j\}) \quad \because (1), \text{ weak union} \quad (2)$$

$$X_k \perp\!\!\!\perp_{\mathcal{P}} X_j \mid \mathbf{X}_{\text{Pre}(k, \pi) \setminus \{j\}} \quad \because (2) \quad (3)$$

where the last formula amounts to  $(j \rightarrow k) \notin E(\mathcal{H}_\pi)$  by (RU).

[ $\Leftarrow$ ] Suppose that  $(j \rightarrow k) \notin E(\mathcal{H}_\pi)$ . We have  $X_k \perp\!\!\!\perp_{\mathcal{P}} X_j \mid \mathbf{X}_{\text{Pre}(k,\pi) \setminus \{j\}}$ . Let  $\mathbf{M}'$  be  $\mathbf{X}_{\text{Pre}(k,\pi) \setminus \{j\}}$ . We have  $X_k \perp\!\!\!\perp_{\mathcal{P}} (\mathbf{X}_{\text{Pre}(k,\pi)} \setminus \mathbf{M}') \mid \mathbf{M}'$  such that  $\mathbf{M}'$  is a Markov blanket of  $X_k$  relative to  $\mathbf{X}_{\text{Pre}(k,\pi)}$ . By **Lemma A.3**,  $X_j \notin \mathbf{M} \subseteq \mathbf{M}'$  and therefore  $(j \rightarrow k) \notin E(\mathcal{G}_\pi)$  by (VP).  $\square$

**Theorem A.5** [Pearl, 1988] *Given a semigraphoid  $\mathcal{P}$  over  $\mathbf{V}$ ,  $\mathcal{G}_\pi$  induced by  $\pi$  using (VP) is Markovian and SGS-minimal for any  $\pi \in \Pi(\mathbf{v})$ .*

**Theorem 3.5** *Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ ,  $\mathcal{G}_\pi$  induced by  $\pi$  using (RU) is Markovian and SGS-minimal for any  $\pi \in \Pi(\mathbf{v})$ .*

*Proof.* Immediate from **Lemma A.4** and **Theorem A.5**.<sup>1</sup>  $\square$

## B CORRECTNESS RESULTS

First, we introduce some permutation-based notations to facilitate our coming proofs. In this section, we use  $\mathcal{G}_\pi$  to denote the DAG induced by  $\pi$  from a graphoid  $\mathcal{P}$  using (RU) unless specified otherwise.

Given a set of variables  $\mathbf{V}$ , consider any  $\pi \in \Pi(\mathbf{v})$  and any pair  $j, k \in \mathbf{v}$  where  $\pi[j] < \pi[k]$ .  $\pi$  can be written as  $\langle \delta_{<j}, j, \delta_{j \sim k}, k, \delta_{>k} \rangle$  such that  $\delta_{<j} = \langle \pi_i : 1 \leq i < \pi[j] \rangle$ ,  $\delta_{j \sim k} = \langle \pi_i : \pi[j] < i < \pi[k] \rangle$ , and  $\delta_{>k} = \langle \pi_i : \pi[k] < i \leq |\pi| \rangle$ . When  $\delta_{j \sim k} = \emptyset$ , we say that  $j$  and  $k$  are  $\pi$ -adjacent. In that case,  $\pi$  can be written as  $\langle \delta_{<j}, j, k, \delta_{>k} \rangle$  instead.

**Definition B.1** *Given a set of variables  $\mathbf{V}$ , for any  $\pi, \tau \in \Pi(\mathbf{v})$ ,*

- (a)  $\tau$  is said to be  $(j, k)$ -different from  $\pi$  for some  $j, k \in \mathbf{v}$  if  $j$  and  $k$  are  $\pi$ -adjacent (i.e.,  $\pi = \langle \delta_{<j}, j, k, \delta_{>k} \rangle$ ) and  $\tau = \langle \delta_{<j}, k, j, \delta_{>k} \rangle$ ;
- (b)  $\pi$  and  $\tau$  are said to be in adjacent transposition (AT) if they are  $(j, k)$ -different for some  $j, k \in \mathbf{v}$ .

**Lemma B.2** *Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , consider any  $\mathcal{H} \in \text{CMC}(\mathcal{P})$ . If  $\pi \in \Pi(\mathbf{v})$  is a causal order of  $\mathcal{G}$ , then  $E(\mathcal{G}_\pi) \subseteq E(\mathcal{H})$ . Also,  $\mathcal{G}_\pi = \mathcal{H}$  if  $\mathcal{H} \in \text{SGS}(\mathcal{P})$ .*

*Proof.* Consider any  $k \in \mathbf{v}$  and  $\text{Nd}(k, \mathcal{H})$  (i.e., the set of  $k$ 's non-descendants in  $\mathcal{H}$ ). Since  $\mathcal{H} \in \text{CMC}(\mathcal{P})$ , it follows that  $X_k \perp\!\!\!\perp_{\mathcal{P}} \mathbf{X}_{\text{Nd}(k, \mathcal{H})} \setminus \mathbf{X}_{\text{Pa}(k, \mathcal{H})} \mid \mathbf{X}_{\text{Pa}(k, \mathcal{H})}$ . Also, we have  $\text{Pa}(k, \mathcal{H}) \subseteq \text{Pre}(k, \pi) \subseteq \text{Nd}(k, \mathcal{H})$  from  $\pi$ 's being a causal order of  $\mathcal{H}$ . By decomposition, we have  $X_k \perp\!\!\!\perp_{\mathcal{P}} \mathbf{X}_{\text{Pre}(k, \pi)} \setminus \mathbf{X}_{\text{Pa}(k, \mathcal{H})} \mid \mathbf{X}_{\text{Pa}(k, \mathcal{H})}$  such that  $\mathbf{X}_{\text{Pa}(k, \mathcal{H})}$  is a Markov blanket of  $X_k$  relative to  $\mathbf{X}_{\text{Pre}(k, \pi)}$ . By **Lemma A.3**, we have  $\text{MB}(X_k, \mathbf{X}_{\text{Pre}(k, \pi)}) \subseteq \text{Pa}(k, \mathcal{H})$ . Consider  $\mathcal{G}_\pi$  induced by (VP). The above entails that  $E(\mathcal{G}_\pi) \subseteq E(\mathcal{H})$  since  $\text{Pa}(k, \mathcal{G}_\pi) = \text{MB}(X_k, \mathbf{X}_{\text{Pre}(k, \pi)}) \subseteq \text{Pa}(k, \mathcal{H})$  for each  $k \in \mathbf{v}$ . Due to **Lemma A.4**,  $E(\mathcal{G}_\pi) \subseteq E(\mathcal{H})$  still holds even if  $\mathcal{G}_\pi$  is induced by (RU). Lastly,  $\mathcal{G}_\pi = \mathcal{H}$  follows from **Definition 3.4** if  $\mathcal{H} \in \text{SGS}(\mathcal{P})$ .  $\square$

**Lemma B.3** [Solus et al., 2021] *Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , consider any  $\pi, \tau \in \Pi(\mathbf{v})$  where  $\tau$  is  $(j, k)$ -different from  $\pi$  for some  $j, k \in \mathbf{v}$ . Then  $\mathcal{G}_\pi = \mathcal{G}_\tau$  if and only if  $X_j \perp\!\!\!\perp_{\mathcal{P}} X_k \mid \mathbf{X}_{\text{Pre}(j, \pi)}$ .*

*Proof.* Suppose that  $X_j \not\perp\!\!\!\perp_{\mathcal{P}} X_k \mid \mathbf{X}_{\text{Pre}(j, \pi)}$ . By (RU), we have  $(j \rightarrow k) \in E(\mathcal{G}_\pi)$ . Note that  $(j \rightarrow k) \notin E(\mathcal{G}_\tau)$  since  $\tau[k] < \tau[j]$  and  $\tau$  is a causal order of  $\mathcal{G}_\tau$  by construction. Hence,  $\mathcal{G}_\pi \neq \mathcal{G}_\tau$ .

On the other hand, suppose that  $X_j \perp\!\!\!\perp_{\mathcal{P}} X_k \mid \mathbf{X}_{\text{Pre}(j, \pi)}$ . Since  $\tau$  is  $(j, k)$ -different from  $\pi$ , we have  $\pi = \langle \delta_{<j}, j, k, \delta_{>k} \rangle$  and  $\tau = \langle \delta_{<j}, k, j, \delta_{>k} \rangle$  according to **Definition B.1** (a). By (RU), we know that  $(k \rightarrow j) \notin E(\mathcal{G}_\tau)$ . Hence,  $\pi$  is a causal order of  $\mathcal{G}_\tau$ . By **Theorem 3.5**,  $\mathcal{G}_\tau \in \text{SGS}(\mathcal{P})$ . Therefore, it follows from **Lemma B.2** that  $\mathcal{G}_\tau = \mathcal{G}_\pi$ .  $\square$

**Lemma B.4** *Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , consider any  $\pi \in \Pi(\mathbf{v})$ . Suppose that  $\mathcal{G}_\pi$  contains a covered edge  $j \rightarrow k$  where  $\pi = \langle \delta_{<j}, j, \delta_{j \sim k}, k, \delta_{>k} \rangle$ . If  $\tau = \langle \delta_{<j}, j, k, \delta_{j \sim k}, \delta_{>k} \rangle$ , then  $\mathcal{G}_\pi = \mathcal{G}_\tau$ .*

*Proof.* Since  $j \rightarrow k$  is a covered edge in  $\mathcal{G}_\pi$ , it follows that  $(i \rightarrow k) \notin E(\mathcal{G}_\pi)$  for each  $i \in \delta_{j \sim k}$ , and thus  $X_i \perp\!\!\!\perp_{\mathcal{P}} X_k \mid \mathbf{X}_{\text{Pre}(i, \pi)}$  by (RU). Hence,  $\mathcal{G}_\pi = \mathcal{G}_\tau$  can be obtained after  $|\delta_{j \sim k}|$  applications of **Lemma B.3**.  $\square$

<sup>1</sup>The two DAG-inducing methods were not differentiated in [Raskutti and Uhler, 2018]. Thus, we provide a proof of **Theorem 3.5**.

**Theorem B.5** [Zhang, 2013] Given a set of variables  $\mathbf{V}$ , for any  $\mathcal{G}, \mathcal{H} \in \text{DAG}(\mathbf{V})$ , if  $\mathbf{E}(\mathcal{G}) \subseteq \mathbf{E}(\mathcal{H})$ , then  $\mathbf{I}(\mathcal{H}) \subseteq \mathbf{I}(\mathcal{G})$ .

**Lemma B.6** [Chickering, 1995] Consider any DAG  $\mathcal{G}$ . Let  $\mathcal{H}$  be the result of reversing  $(i \rightarrow j) \in \mathbf{E}(\mathcal{G})$ . Then  $\mathcal{H} \in \text{MEC}(\mathcal{G})$  if and only if  $i \rightarrow j$  is a covered edge.

**Theorem B.7** [Chickering, 1995] Consider any pair of DAGs  $\mathcal{G}$  and  $\mathcal{H}$  over the same set of variables s.t.  $\mathcal{H} \in \text{MEC}(\mathcal{G})$ , and for which there are  $k$  edges in  $\mathcal{G}$  that have opposite orientation in  $\mathcal{H}$ . Then there exists a sequence of  $k$  distinct covered edge reversals in  $\mathcal{G}$  s.t.  $\mathcal{G}$  becomes  $\mathcal{H}$  after all reversals.

**Lemma B.8** Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , consider any  $\pi \in \Pi(\mathbf{v})$ . Suppose that  $(j \rightarrow k) \in \mathbf{E}(\mathcal{G}_\pi)$  is a covered edge, and let  $\mathcal{H}$  be the DAG resulted from reversing  $(j \rightarrow k)$  in  $\mathcal{G}_\pi$ . If  $\tau = \text{tuck}(\pi, j, k)$ , then

- (a)  $\tau$  is a causal order of  $\mathcal{H}$ ;
- (b)  $\mathbf{E}(\mathcal{G}_\tau) \subseteq \mathbf{E}(\mathcal{H})$ ;
- (c)  $|\mathbf{E}(\mathcal{G}_\tau)| \leq |\mathbf{E}(\mathcal{G}_\pi)|$ ;
- (d)  $\mathbf{I}(\mathcal{G}_\pi) \subseteq \mathbf{I}(\mathcal{G}_\tau)$ .

*Proof.* A similar lemma has been shown in [Solus et al., 2021]. First, we write  $\pi = \langle \delta_{<j, j}, \delta_{j \sim k}, k, \delta_{>k} \rangle$  as usual. Consider  $\pi' = \langle \delta_{<j, j}, k, \delta_{j \sim k}, \delta_{>k} \rangle$ . By **Lemma B.4**, we have  $\mathcal{G}_\pi = \mathcal{G}_{\pi'}$ . Note that  $\tau = \text{tuck}(\pi, j, k) = \langle \delta_{<j, k}, j, \delta_{j \sim k}, \delta_{>k} \rangle$  because  $\text{tuck}(\pi, j, k)$  is a covered tuck. Thus,  $\tau$  is  $(j, k)$ -different from  $\pi'$ . Also, since  $\pi'$  is a causal order of  $\mathcal{G}_\pi$ , it follows that  $\tau$  is a causal order of  $\mathcal{H}$  and thus (a) is proven.

Next, observe that  $\mathbf{I}(\mathcal{G}_\pi) = \mathbf{I}(\mathcal{H})$  from **Lemma B.6**. From  $\mathcal{G}_\pi \in \text{CMC}(\mathcal{P})$  by **Theorem 3.5**, we know that  $\mathcal{H} \in \text{CMC}(\mathcal{P})$ . Thus, (b) immediately follows from (a) and **Lemma B.2**. Also, (c) is entailed by  $|\mathbf{E}(\mathcal{G}_\tau)| \leq |\mathbf{E}(\mathcal{H})| = |\mathbf{E}(\mathcal{G}_\pi)|$ . Finally, by **Theorem B.5**, we have  $\mathbf{I}(\mathcal{G}_\pi) = \mathbf{I}(\mathcal{H}) \subseteq \mathbf{I}(\mathcal{G}_\tau)$  as desired in (d).  $\square$

Before we compare TSP and unbounded GRaSP<sub>0</sub>, we want to make an assumption related to how the set of covered edges in any particular DAG is ordered. To see the importance of such an assumption, observe that different orderings of  $\mathbf{E}^0(\mathcal{G}_\pi)$  (i.e., the set of covered edges in an induced DAG  $\mathcal{G}_\pi$ ) can alter the output of TSP and also GRaSP<sub>0</sub>. For example, suppose that  $(j \rightarrow k), (j' \rightarrow k') \in \mathbf{E}^0(\mathcal{G}_\pi)$ . Say the DFS of GRaSP<sub>0</sub> starts with performing  $\text{tuck}(\pi, j, k)$  and leads to some permutation  $\tau$ . However, choosing to perform  $\text{tuck}(\pi, j', k')$  instead at the beginning of the DFS procedure can lead to some  $\tau'$  where  $\mathcal{G}_\tau \neq \mathcal{G}_{\tau'}$ . Hence, we enforce the assumption that the ordering of  $\mathbf{E}^0(\mathcal{G})$  for any DAG  $\mathcal{G}$  is fixed arbitrarily. For instance,  $(j \rightarrow k)$  precedes  $(j' \rightarrow k')$  in  $\mathbf{E}^0(\mathcal{G})$  if  $j < j'$ , or  $j = j'$  and  $k < k'$ . Consequently, the issue of order-dependence can be avoided even when comparing a Chickering sequence found by TSP and a ct-sequence found by unbounded GRaSP<sub>0</sub>. In the following, this assumption will be made implicitly.

Now we revisit how TSP works. Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$  and an initial permutation  $\pi \in \Pi(\mathbf{v})$ , TSP begins with setting  $\mathcal{G}$  as the induced  $\mathcal{G}_\pi$ . Starting with the root  $\mathcal{G}$ , TSP performs DFS to identify a SGS-minimal DAG  $\mathcal{H}$  connected by a Chickering sequence from  $\mathcal{G}$  such that  $|\mathbf{E}(\mathcal{G})| > |\mathbf{E}(\mathcal{H})|$ . TSP returns  $\mathcal{G}$  if no such  $\mathcal{H}$  is found. Otherwise, it updates  $\mathcal{G}$  as  $\mathcal{H}$  and repeat the procedure.

The DFS procedure of TSP aims to traverse from one SGS-minimal DAG to another SGS-minimal DAG by the construction of a Chickering sequence. Though we know that a Chickering sequence is obtained by the reversals of covered edges and deletions of directed edges, Solus et al. [2021] did not specify any ordering of these operations. Below we provide a more precise definition of the Chickering sequences considered by TSP.

**Definition B.9** Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , a TSP-Chickering sequence  $\mathfrak{C} = \langle \mathcal{G}^1, \dots, \mathcal{G}^m \rangle$  is a Chickering sequence satisfying the following condition:

- (a)  $\mathcal{G}^1, \mathcal{G}^m \in \text{SGS}(\mathcal{P})$ ;
- (b)  $\mathcal{G}^i$  and  $\mathcal{G}^{i'}$  are pairwise distinct for  $1 \leq i < i' \leq m$ ;
- (c) if  $|\mathbf{E}(\mathcal{G}^1)| = |\mathbf{E}(\mathcal{G}^m)|$ , then  $\mathcal{G}^1, \dots, \mathcal{G}^m \in \text{SGS}(\mathcal{P})$  where they differ by the reversals of some covered edges;
- (d) otherwise, there exists a turning index  $1 < l < m$  such that (i)  $\mathcal{G}^1, \dots, \mathcal{G}^{l-1} \in \text{SGS}(\mathcal{P})$ , (ii)  $\mathcal{G}^1, \dots, \mathcal{G}^l$  differ by the reversals of some covered edges, and (iii)  $\mathcal{G}^{i+1}$  is obtained from deleting a directed edge in  $\mathcal{G}^i \notin \text{SGS}(\mathcal{P})$  for each  $l \leq i < m$ .

Readers are suggested to find the original pseudocode of TSP in Solus et al. [2021] to verify that our **Definition B.9** is a fair description of the Chickering sequences considered by TSP. Conditions (a) and (b) are straightforward. (c) refers to the case where TSP cannot find a sparser SGS-minimal DAG. So if any  $\mathcal{G}^i$  in  $\mathfrak{C}$  were non-SGS-minimal, then TSP would have obtained a proper subgraph of  $\mathcal{G}^i$  which is SGS-minimal by a series of edge-deletion. (d) refers to the case where TSP manages to find a sparser SGS-minimal DAG. Notice that  $\mathcal{G}^2$  must be obtained by a covered edge reversal from  $\mathcal{G}^1$  since  $\mathcal{G}^1 \in \text{SGS}(\mathcal{P})$ . If  $\mathcal{G}^2 \notin \text{SGS}(\mathcal{P})$ , then TSP can obtain the desired SGS-minimal DAG by a series of edge-deletion from  $\mathcal{G}^2$ . But if  $\mathcal{G}^2 \in \text{SGS}(\mathcal{P})$ , the procedure above repeats until finding the turning index  $l$  such that  $\mathcal{G}^l \notin \text{SGS}(\mathcal{P})$  and then the sparser  $\mathcal{G}^m \in \text{SGS}(\mathcal{P})$  can be obtained by a series of edge-deletion from  $\mathcal{G}^l$ .

Now we compare TSP and unbounded GRaSP<sub>0</sub> by considering their respective sequences traversed in the DFS procedure.

**Lemma B.10** *Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , consider any  $\pi \in \Pi(\mathbf{v})$  and  $\tau = \text{tuck}(\pi, j, k)$  where  $(j \rightarrow k) \in \mathbf{E}^0(\mathcal{G}_\pi)$ . Given that  $\mathfrak{T} = \langle \pi, \tau \rangle$  is a ct-sequence,*

- (a) *if  $|\mathbf{E}(\mathcal{G}_\tau)| = |\mathbf{E}(\mathcal{G}_\pi)|$ , then  $\mathfrak{C} = \langle \mathcal{G}_\pi, \mathcal{G}_\tau \rangle$  is a TSP-Chickering sequence where  $\mathcal{G}_\tau$  is obtained from reversing  $(j \rightarrow k) \in \mathbf{E}^0(\mathcal{G}_\pi)$ ;*
- (b) *otherwise, there exists a TSP-Chickering sequence  $\mathfrak{C} = \langle \mathcal{G}_\pi = \mathcal{G}^1, \dots, \mathcal{G}^m = \mathcal{G}_\tau \rangle$  s.t.  $\mathcal{G}^2$  is obtained from reversing  $(j \rightarrow k) \in \mathbf{E}^0(\mathcal{G}_\pi)$ , and  $\mathcal{G}^{i+1}$  is obtained from deleting a directed edge in  $\mathcal{G}^i$  for each  $2 \leq i < m$ .*

*Proof.* First, consider the DAG  $\mathcal{H}$  obtained from reversing  $(j \rightarrow k) \in \mathbf{E}^0(\mathcal{G}_\pi)$ . We start with the case in (a) where  $|\mathbf{E}(\mathcal{G}_\tau)| = |\mathbf{E}(\mathcal{G}_\pi)| = |\mathbf{E}(\mathcal{H})|$ . We want to show that  $\mathcal{G}_\tau = \mathcal{H}$ . By **Lemma B.8** (b), we have  $\mathbf{E}(\mathcal{G}_\tau) \subseteq \mathbf{E}(\mathcal{H})$ . If  $\mathbf{E}(\mathcal{G}_\tau) \subset \mathbf{E}(\mathcal{H})$  holds, then  $|\mathbf{E}(\mathcal{G}_\tau)| = |\mathbf{E}(\mathcal{H})|$  will be violated. Hence, we have  $\mathcal{G}_\tau = \mathcal{H}$  and thus  $\mathfrak{C} = \langle \mathcal{G}_\pi, \mathcal{H} = \mathcal{G}_\tau \rangle$  is our desired TSP-Chickering sequence.

For (b), it follows from **Lemma B.8** (c) that  $|\mathbf{E}(\mathcal{G}_\tau)| < |\mathbf{E}(\mathcal{H})| = |\mathbf{E}(\mathcal{G}_\pi)|$ . Let  $\mathcal{G}_\pi$  and  $\mathcal{H}$  be  $\mathcal{G}^1$  and  $\mathcal{G}^2$  respectively. By **Lemma B.8** (b) again, we have  $\mathbf{E}(\mathcal{G}_\tau) \subset \mathbf{E}(\mathcal{G}^2)$  such that we can remove a directed edge from  $\mathcal{G}^2$  once at a time until obtaining  $\mathcal{G}_\tau$ . Therefore, we have the desired TSP-Chickering sequence in (b).  $\square$

**Lemma B.11** *Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , consider any TSP-Chickering sequence  $\mathfrak{C} = \langle \mathcal{G}^1, \dots, \mathcal{G}^m \rangle$ . Let  $\pi^1$  be a causal order of  $\mathcal{G}^1$ . Then*

- (a) *if  $|\mathbf{E}(\mathcal{G}^1)| = |\mathbf{E}(\mathcal{G}^m)|$ , then  $\mathcal{G}^{i+1} = \mathcal{G}_{\pi^{i+1}} = \mathcal{G}_{\text{tuck}(\pi, j, k)}$  where  $j \rightarrow k$  is the covered edge reversed to obtain  $\mathcal{G}^{i+1}$  from  $\mathcal{G}^i$  for each  $1 \leq i < m$  s.t.  $\mathfrak{T} = \langle \pi^1, \dots, \pi^m \rangle$  is a ct-sequence;*
- (b) *otherwise, then  $\mathcal{G}^{i+1} = \mathcal{G}_{\pi^{i+1}} = \mathcal{G}_{\text{tuck}(\pi, j, k)}$  where  $j \rightarrow k$  is the covered edge reversed to obtain  $\mathcal{G}^{i+1}$  from  $\mathcal{G}^i$  for each  $1 \leq i < l$  where  $l$  is the turning index of  $\mathfrak{C}$  and  $\mathcal{G}_{\pi^l} = \mathcal{G}^m$  s.t.  $\mathfrak{T} = \langle \pi^1, \dots, \pi^l \rangle$  is a ct-sequence.*

*Proof.* (a) can be easily shown by **Lemma B.8**(a) and **Lemma B.2**. For (b), the proof of  $\mathcal{G}^i = \mathcal{G}_{\pi^i}$  for each  $1 \leq i < l$  is similar to that in (a). So we consider  $l$  where  $\mathcal{G}^l \notin \text{SGS}(\mathcal{P})$  according to **Definition B.9**(d). However, it follows from **Lemma B.8** (a) that  $\pi^l$  is a causal order of  $\mathcal{G}^l$ . Since  $\mathbf{E}(\mathcal{G}^m) \subset \mathbf{E}(\mathcal{G}^l)$ , we know that  $\pi^l$  is also a causal order of  $\mathcal{G}^m$ . Lastly, given that  $\mathcal{G}^m \in \text{SGS}(\mathcal{P})$ , it follows from **Lemma B.2** that  $\mathcal{G}_{\pi^l} = \mathcal{G}^m$ .  $\square$

**Lemma 4.4** *Given a graphoid  $\mathcal{P}$ , for any  $\pi \in \Pi(\mathbf{v})$  and any Chickering sequence from  $\mathcal{G}_\pi$  to some  $\mathcal{H} \in \text{SGS}(\mathcal{P})$  considered by TSP, there exists a ct-sequence  $\langle \pi, \dots, \tau \rangle$  s.t.  $\mathcal{G}_\tau = \mathcal{H}$ .*

*Proof.* Given that a Chickering sequence considered by TSP is simply a TSP-Chickering sequence defined in **Definition B.9**, the lemma follows immediately from **Lemma B.11**.  $\square$

**Theorem 4.7** *Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$  and any initial permutation  $\pi \in \Pi(\mathbf{v})$ , the DAG induced by the output of unbounded GRaSP<sub>0</sub> is equivalent to the DAG returned by TSP.*

*Proof.* Immediate from **Lemma B.10** and **Lemma B.11**.  $\square$

Now we turn to the discussion on the correctness of GRaSP<sub>0</sub> under faithfulness.

**Lemma B.12** Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$  and any  $\pi \in \Pi(\mathbf{v})$ , if  $\mathcal{G}_\pi \notin \text{Pm}(\mathcal{P})$ , then there exists a ct-sequence  $\mathfrak{T} = \langle \pi, \dots, \tau \rangle$  s.t.  $\text{I}(\mathcal{G}_\pi) \subset \text{I}(\mathcal{G}_\tau)$ .

*Proof.* Suppose that  $\mathcal{G}_\pi \notin \text{Pm}(\mathcal{P})$ . By **Definition 3.7**, it follows that there exists  $\mathcal{H} \in \text{CMC}(\mathcal{P})$  s.t.  $\text{I}(\mathcal{G}_\pi) \subset \text{I}(\mathcal{H}) \subseteq \text{I}(\mathcal{P})$ . By **Theorem 3.6**, we know that there exists a Chickering sequence  $\mathfrak{C}_0 = \langle \mathcal{G}_\pi = \mathcal{G}^1, \dots, \mathcal{G}^l = \mathcal{H} \rangle$ . Without loss of generality, suppose that  $\mathfrak{C}_0$  is the shortest Chickering sequence where each  $\mathcal{G}^{i+1}$  differs from  $\mathcal{G}^i \in \text{SGS}(\mathcal{P})$  by the reversal of a covered edge in  $\text{E}^0(\mathcal{G}^i)$  for each  $1 \leq i < l - 1$ , and  $\mathcal{G}^l$  is obtained from deleting a directed edge in  $\mathcal{G}^{l-1}$ . Notice that  $|\text{E}(\mathcal{G}^l)| < |\text{E}(\mathcal{G}^1)|$  due to the edge deletion. If  $\mathcal{G}^l \in \text{SGS}(\mathcal{P})$ , then  $\mathfrak{C}_0$  is a TSP-Chickering sequence. Otherwise, we can easily construct a TSP-Chickering sequence  $\mathfrak{C} = \langle \mathcal{G}^1, \dots, \mathcal{G}^m \rangle$  with  $l - 1$  as the turning index and  $\mathcal{G}^m \in \text{SGS}(\mathcal{P})$  obtained by repeated edge-deletion from  $\mathcal{G}^l$  such that  $\text{I}(\mathcal{G}^1) \subset \text{I}(\mathcal{G}^l) \subset \text{I}(\mathcal{G}^m)$ . By **Lemma B.11** (b), we have the desired ct-sequence.  $\square$

**Theorem 4.5** Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$  and any  $\pi \in \Pi(\mathbf{v})$ , if  $\mathcal{G}_\pi \notin \text{Pm}(\mathcal{P})$ , then there exists a ct-sequence  $\mathfrak{T} = \langle \pi, \dots, \tau \rangle$  s.t.  $\mathcal{G}_\tau \in \text{Pm}(\mathcal{P})$ .

*Proof.* Immediate from **Lemma B.12**.  $\square$

**Theorem B.13** Unbounded  $\text{GRaSP}_0$  is correct and pointwise consistent under faithfulness.

*Proof.* We review the argument for the correctness of unbounded  $\text{GRaSP}_0$  under faithfulness given in the main paper. Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , consider any initial permutation  $\pi \in \Pi(\mathbf{v})$ . Given that unbounded  $\text{GRaSP}_0$  greedily search for a ct-sequence from  $\pi$ , it is guaranteed by **Theorem 4.5** that  $\tau$  returned by unbounded  $\text{GRaSP}_0$  in **Algorithm 2** induces a P-minimal DAG. Under faithfulness, we have  $\mathcal{G}_\tau \in \text{MEC}(\mathcal{G}^*)$  due to **Theorem 3.8** and hence unbounded  $\text{GRaSP}_0$  is correct.

Alternatively, the correctness and pointwise consistency of unbounded  $\text{GRaSP}_0$  can also be proven directly from **Theorem 4.7** and the corresponding results of TSP in [Solus et al., 2021].  $\square$

**Corollary 4.6** Unbounded  $\text{GRaSP}_0$ ,  $\text{GRaSP}_1$ , and  $\text{GRaSP}_2$  are correct and pointwise consistent under faithfulness.

In the following, we want to prove that faithfulness is not only sufficient, but also *necessary* for the correctness of TSP and unbounded  $\text{GRaSP}_0$ . We first want to prove an interesting and novel equivalence between two causal razors: faithfulness and u-P-minimality.

**Lemma B.14** Given a joint probability distribution  $\mathcal{P}$  over  $\mathbf{V}$ , for any  $\langle X_i, X_j \mid \mathbf{X}_k \rangle \in \text{I}(\mathcal{P})$ , there exists  $\mathcal{G} \in \text{DAG}(\mathbf{V})$  s.t.  $\text{I}(\mathcal{G}) = \{\langle X_i, X_j \mid \mathbf{X}_k \rangle\}$ .

*Proof.* Consider  $\mathbf{V} = \{X_1, \dots, X_m\}$ . An empty DAG suffices when  $m = 2$ . So assume that  $m \geq 3$ . Without loss of generality, consider  $\langle X_1, X_{k+2} \mid \mathbf{X}_k \rangle \in \text{I}(\mathcal{P})$  where  $\mathbf{k} = \langle 2, \dots, k + 1 \rangle$ , and the remaining vertices are  $\langle k + 3, \dots, m \rangle$ . We propose a procedure which guarantees the existence of the desired DAG  $\mathcal{G}$ .

- 
- 1  $\mathcal{G} \leftarrow$  a complete undirected graph over  $\mathbf{v}$
  - 2 remove the adjacency  $1 - k + 2$  in  $\mathcal{G}$
  - 3 **foreach**  $(j, k)$  that are adjacent in  $\mathcal{G}$  **do**
  - 4     **if**  $j < k$  **then**
  - 5         orient  $j \rightarrow k$  in  $\mathcal{G}$
  - 6 **return**  $\mathcal{G}$
- 

Line 3 to 5 guarantee that  $\mathcal{G}$  is a DAG since all edges are directed and pointing from lower indices to higher indices such that no directed cycle can occur. Finally,  $1 \perp_{\mathcal{G}} k + 2 \mid \mathbf{k}$  holds because all paths from 1 to  $k + 2$  either contain a non-collider  $i \in \mathbf{k}$  or contain a collider  $i \notin \mathbf{k}$ . Therefore,  $\text{I}(\mathcal{G}) = \{\langle X_1, X_{k+2} \mid \mathbf{X}_k \rangle\}$  because no other d-separation relations hold in  $\mathcal{G}$ .  $\square$

**Theorem B.15** For any joint probability distribution  $\mathcal{P}$ ,  $\text{CFC}(\mathcal{P}) = \text{uPm}(\mathcal{P})$ .

*Proof.*  $\subseteq$  Suppose that  $\mathcal{G} \in \text{CFC}(\mathcal{P})$ . It follows that  $\mathcal{G} \in \text{Pm}(\mathcal{P})$  by **Definition 3.7**. For any  $\mathcal{G}' \in \text{CMC}(\mathcal{P})$ , if  $\text{I}(\mathcal{G}') \subset \text{I}(\mathcal{G})$ , then  $\mathcal{G}' \notin \text{Pm}(\mathcal{P})$ . Hence, if  $\mathcal{G}' \in \text{Pm}(\mathcal{P})$ , then  $\text{I}(\mathcal{G}') = \text{I}(\mathcal{G})$ . Hence,  $\mathcal{G} \in \text{uPm}(\mathcal{P})$ .

$\supseteq$  Suppose that  $\mathcal{G} \notin \text{CFC}(\mathcal{P})$ . Since  $\text{uPm}(\mathcal{P}) \subseteq \text{Pm}(\mathcal{P})$  by **Definition 3.7**, if  $\mathcal{G} \notin \text{Pm}(\mathcal{P})$ , we have  $\mathcal{G} \notin \text{uPm}(\mathcal{P})$  immediately. So consider the case where  $\mathcal{G} \in \text{Pm}(\mathcal{P})$ . It follows from  $\mathcal{G} \notin \text{CFC}(\mathcal{P})$  that there exists a CI relation  $\psi \in \text{I}(\mathcal{P}) \setminus \text{I}(\mathcal{G})$ . By **Lemma B.14**, we can construct a DAG  $\mathcal{G}^0$  such that  $\text{I}(\mathcal{G}^0) = \{\psi\}$ . Consequently, there exists  $\mathcal{G}^1 \in \text{Pm}(\mathcal{P})$  such that  $\text{I}(\mathcal{G}^0) \subseteq \text{I}(\mathcal{G}^1) \subseteq \text{I}(\mathcal{P})$ . Since  $\psi \in \text{I}(\mathcal{G}^1)$ , we know that  $\mathcal{G}^1 \notin \text{MEC}(\mathcal{G})$ . Given that both  $\mathcal{G}, \mathcal{G}^1 \in \text{Pm}(\mathcal{P})$ , we have  $\mathcal{G} \notin \text{uPm}(\mathcal{P})$ .  $\square$

**Theorem 4.8** Given a graphoid  $\mathcal{P}$ , faithfulness is necessary for the correctness of TSP.

*Proof.* Suppose that  $(\mathcal{G}^*, \mathcal{P})$  is unfaithful. We consider the two kinds of unfaithfulness in [Zhang and Spirtes, 2008]: *detectable* (i.e.,  $\text{CFC}(\mathcal{P}) = \emptyset$ ) versus *undetectable* (i.e.,  $\mathcal{G}' \in \text{CFC}(\mathcal{P})$  where  $\mathcal{G}' \notin \text{MEC}(\mathcal{G}^*)$ ). For the latter, TSP can identify  $\mathcal{G}_\tau \in \text{Pm}(\mathcal{P}) = \text{CFC}(\mathcal{P}) = \text{MEC}(\mathcal{G}')$ . However, TSP is incorrect because  $\mathcal{G}_\tau \notin \text{MEC}(\mathcal{G}^*)$ .

On the other hand, consider the case that  $\text{CFC}(\mathcal{P}) = \emptyset$ . By **Theorem B.15**, there exists  $\mathcal{G} \in \text{Pm}(\mathcal{P})$  such that  $\mathcal{G} \notin \text{MEC}(\mathcal{G}^*)$  even if  $\mathcal{G}^* \in \text{Pm}(\mathcal{P})$ . Recall that Chickering algorithm can only allow us to traverse to a DAG  $\mathcal{H}$  from  $\mathcal{G}$  satisfying  $\text{I}(\mathcal{G}) \subseteq \text{I}(\mathcal{H})$ . It entails that Chickering algorithm can only obtain DAGs that are in  $\text{MEC}(\mathcal{G})$  since  $\mathcal{G} \in \text{Pm}(\mathcal{P})$  and hence never be able to reach  $\mathcal{G}^*$  where  $\text{I}(\mathcal{G}_\pi) \not\subseteq \text{I}(\mathcal{G}^*)$ . Therefore, by setting  $\pi$  as the initial permutation to TSP where  $\mathcal{G}_\pi = \mathcal{G}$ , TSP will return  $\mathcal{G}_\pi$  incorrectly.  $\square$

Notice that **Theorem 4.8** is contrary to what Solus et al. [2021] suggested. They proposed an example arguing that TSP can be correct even under (detectable) unfaithfulness.<sup>2</sup> However, the distribution used in the example is not a semigraphoid. This renders their example illegitimate because every joint probability distribution is a semigraphoid.

## C ESP AND GRASP-1

As shown in **Theorem 4.8** in the last section, TSP cannot be correct under unfaithfulness by choosing an arbitrary initial permutation. Consequently, one important question is how to relax the search space of TSP to identify a sparser permutation under unfaithfulness. Solus et al. [2021] proposed the *Edge SP* (ESP) algorithm based on an assumption strictly weaker than that assumed by TSP. However, unlike TSP, they did not provide an operational version of ESP in their work. In this section, we are going to show a theorem similar to **Theorem 4.7** but with respect to ESP and unbounded GRASP<sub>1</sub>. In other words, unbounded GRASP<sub>1</sub> is an operational version of ESP. In the following, we first examine some technical notations used in [Mohammadi et al., 2018] and [Solus et al., 2021]. Readers are strongly suggested to visit [Solus et al., 2021] for the full discussion of ESP and relevant notations.

Given a set of measured variables  $\mathbf{V}$ , a *permutohedron* on  $\mathbf{v}$ , denoted  $A_{\mathbf{v}}$ , is the convex hull in  $\mathbb{R}^{|\mathbf{v}|}$  of all permutations in  $\Pi(\mathbf{v})$ . In simpler terms,  $A_{\mathbf{v}}$  is the *state space* with each *state* being a permutation  $\pi \in \Pi(\mathbf{v})$ . The neighborhood of states in  $A_{\mathbf{v}}$  is defined by adjacent transpositions (ATs) as in **Definition B.1** (b).

Notice that different states in  $A_{\mathbf{v}}$  can induce the same DAG given a graphoid  $\mathcal{P}$ . Thus, a natural way to narrow down the search space is to identify permutations inducing the same DAG. **Lemma B.3** provides such a characterization. Construct  $A_{\mathbf{v}}(\mathcal{P})$  by *contracting* neighborhood in  $A_{\mathbf{v}}$  to ATs that correspond to the CI relations in  $\text{I}(\mathcal{P})$  specified in **Lemma B.3**. To be more specified, the contracted permutohedron  $A_{\mathbf{v}}(\mathcal{P})$ , also known as the *DAG associahedron*, is the state space with each state being an induced DAG.<sup>3</sup> Two states  $\mathcal{G}_1, \mathcal{G}_2$  in  $A_{\mathbf{v}}(\mathcal{P})$  are neighbors if and only if there exist  $\pi^1, \pi^2 \in \Pi(\mathbf{v})$  s.t.  $\mathcal{G}_{\pi^1} = \mathcal{G}_1, \mathcal{G}_{\pi^2} = \mathcal{G}_2$ , and  $\pi^1$  and  $\pi^2$  are neighbors in the permutohedron  $A_{\mathbf{v}}$ . As shown by Mohammadi et al. [2018], the DAG associahedron is a convex polytope where each vertex of  $A_{\mathbf{v}}(\mathcal{P})$  corresponds to a different DAG.

To draw a clearer picture, consider any  $\pi, \tau \in \Pi(\mathbf{v})$  where  $\tau$  is  $(j, k)$ -different from  $\pi$  for some  $j, k \in \mathbf{v}$ . They are neighbors in  $A_{\mathbf{v}}$  but they do not necessarily induce the same DAG. If  $X_j \perp\!\!\!\perp_{\mathcal{P}} X_k \mid \mathbf{X}_{\text{Pre}(j, \pi)}$  holds, they induce the same DAG and

<sup>2</sup>See Figure 2 in the [supplementary materials](#) of [Solus et al., 2021].

<sup>3</sup>One can equivalently express each state in the DAG associahedron as the set of permutations which induce the same DAG. This is the original representation in [Mohammadi et al., 2018]. However, we prefer the representation given in [Solus et al., 2021] in the sense that one can easily compare DAGs that are in neighborhood.

thus correspond to the same state  $\mathcal{G}_\pi$  in the DAG associahedron  $A_{\mathbf{v}}(\mathcal{P})$ . But if the CI relation does not hold, then  $\mathcal{G}_\pi$  and  $\mathcal{G}_\tau$  are neighbors in  $A_{\mathbf{v}}(\mathcal{P})$ . See Figure 1 for an example from [Solus et al., 2021].

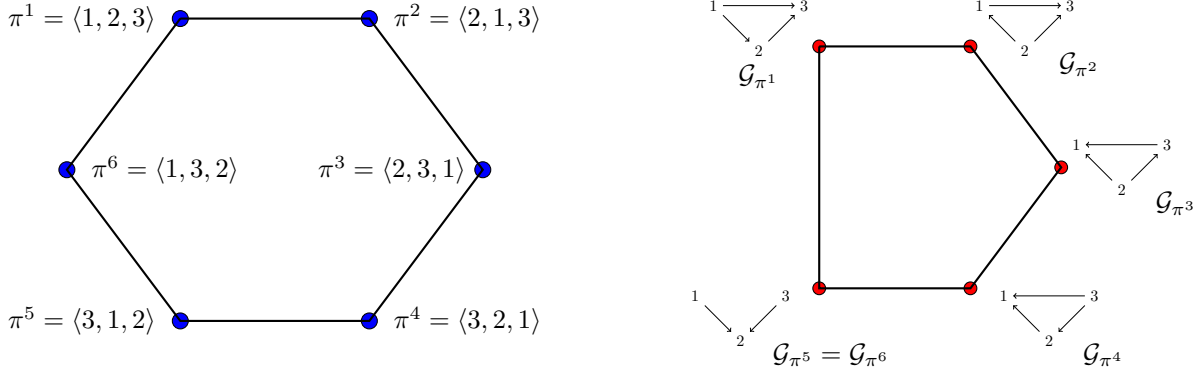


Figure 1: Given  $\mathbf{V} = \{X_1, X_2, X_3\}$ , consider  $\mathcal{I}(\mathcal{P}) = \{\langle X_1, X_3 \mid \emptyset \rangle\}$ . The diagram on the left is the permutohedron  $A_{\mathbf{v}}$  where each state is a permutation in  $\Pi(\mathbf{v})$ . The one on the right is the DAG associahedron  $A_{\mathbf{v}}(\mathcal{P})$  where each state is a different DAG in  $\text{SGS}(\mathcal{P})$ . In particular, the two states  $\pi^5$  and  $\pi^6$  in  $A_{\mathbf{v}}$  are collapsed into a single state in  $A_{\mathbf{v}}(\mathcal{P})$  because they induce the same DAG.

Observe that each state in the DAG associahedron  $A_{\mathbf{v}}(\mathcal{P})$  corresponds to a SGS-minimal DAG according to **Theorem 3.5**. ESP performs a greedy DFS in  $A_{\mathbf{v}}(\mathcal{P})$ . Given an initial permutation  $\pi \in \Pi(\mathbf{v})$ , set  $\mathcal{G}$  as the induced  $\mathcal{G}_\pi$  and traverse through  $A_{\mathbf{v}}(\mathcal{P})$  by a *weakly decreasing walk* to obtain  $\mathcal{H}$  where  $|\mathcal{E}(\mathcal{H})| < |\mathcal{E}(\mathcal{G}_\pi)|$ .<sup>4</sup> If no such  $\mathcal{H}$  exists, ESP returns  $\mathcal{G} = \mathcal{G}_\pi$ ; else  $\mathcal{G}$  is reset as  $\mathcal{H}$  and repeat.

As noted by Solus et al. [2021], the construction of  $A_{\mathbf{v}}(\mathcal{P})$  is inefficient since one is only required to know the neighboring states instead of the entire  $A_{\mathbf{v}}(\mathcal{P})$  to perform the traversal. Below we show that unbounded  $\text{GRaSP}_1$  can efficiently learn the neighbors of each state in  $A_{\mathbf{v}}(\mathcal{P})$  by our permutation-based operation *tuck* performed on singular edges. Before examining this claim, we introduce some useful definitions.

Given the permutohedron  $A_{\mathbf{v}}$ , a *walk*  $\mathfrak{W} = \langle \pi^1, \dots, \pi^m \rangle$  is a sequence of neighboring states in  $A_{\mathbf{v}}$  such that  $\pi^i, \pi^{i+1} \in \Pi(\mathbf{v})$  are in AT for each  $1 \leq i < m$ .

**Definition C.1** Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , for any walk  $\mathfrak{W} = \langle \pi^1, \dots, \pi^m \rangle$  in  $A_{\mathbf{v}}$ ,

- (a)  $\mathfrak{W}$  is said to be *DAG-preserving* if  $\mathcal{G}_{\pi^1} = \dots = \mathcal{G}_{\pi^m}$ ;
- (b)  $\mathfrak{W}$  is said to be *DAG-changing* if  $\langle \pi^1, \dots, \pi^{m-1} \rangle$  is DAG-preserving and  $\mathcal{G}_{\pi^{m-1}} \neq \mathcal{G}_{\pi^m}$ .

In addition, for each DAG-changing walk  $\mathfrak{W} = \langle \pi^1, \dots, \pi^m \rangle$ , we say that  $\mathfrak{W}$  is relative to  $(j, k)$  if  $\pi^m$  is  $(j, k)$ -different from  $\pi^{m-1}$  for some  $j, k \in \mathbf{v}$ . Thus, each DAG-changing walk is relative to a pair of vertices corresponding to the last AT performed in the walk.

**Definition C.2** Given a set of variables  $\mathbf{V}$ , consider any  $j, k \in \mathbf{v}$ . Two DAGs  $\mathcal{G}, \mathcal{H} \in \text{DAG}(\mathbf{V})$  are said to be  $(j, k)$ -reverse if  $(j \rightarrow k) \in \mathcal{E}(\mathcal{G})$  and  $(k \rightarrow j) \in \mathcal{E}(\mathcal{H})$ , and there does not exist any other  $j', k' \in \mathbf{v}$  s.t.  $(j' \rightarrow k') \in \mathcal{E}(\mathcal{G})$  and  $(k' \rightarrow j') \in \mathcal{E}(\mathcal{H})$ .

**Lemma C.3** Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , consider any  $j, k \in \mathbf{v}$ . Then there exists a DAG-changing walk  $\mathfrak{W} = \langle \pi^1, \dots, \pi^m \rangle$  relative to  $(j, k)$  in  $A_{\mathbf{v}}$  if and only if  $\mathcal{G}_{\pi^1}$  and  $\mathcal{G}_{\pi^m}$  are neighbors in  $A_{\mathbf{v}}(\mathcal{P})$  that are  $(j, k)$ -reverse.

*Proof.* For the forward direction, given that  $\pi^{m-1}$  and  $\pi^m$  are  $(j, k)$ -different but induce different DAGs, it follows from the definition of  $A_{\mathbf{v}}(\mathcal{P})$  that  $\mathcal{G}_{\pi^1} = \mathcal{G}_{\pi^{m-1}}$  and  $\mathcal{G}_{\pi^m}$  are neighbors in  $A_{\mathbf{v}}(\mathcal{P})$ . Also, we know that  $(j \rightarrow k) \in \mathcal{E}(\mathcal{G}_{\pi^{m-1}})$  and

<sup>4</sup>In [Solus et al., 2021], their pseudocode does not indicate that such a walk needs to be weakly decreasing but such a requirement is imposed in the description of the algorithm.



$(k \rightarrow j) \in \mathbf{E}(\mathcal{G}_{\pi^m})$  by **Lemma B.3** and (RU). The fact that  $\mathcal{G}_{\pi^1} = \mathcal{G}_{\pi^{m-1}}$  and  $\mathcal{G}_{\pi^m}$  are  $(j, k)$ -reverse follows immediately from (RU) and the assumption that  $\pi^{m-1}$  is  $(j, k)$ -different from  $\pi^m$ .

For the backward direction, suppose that  $\mathcal{G}_\pi$  and  $\mathcal{G}_\tau$  are neighbors in  $\mathbf{A}_\mathbf{v}(\mathcal{P})$  that are  $(j, k)$ -reverse. It entails from (RU) that there exist  $\pi', \tau' \in \Pi(\mathbf{v})$  such that  $\pi'$  and  $\tau'$  are  $(j, k)$ -different where  $\mathcal{G}_\tau = \mathcal{G}_{\tau'}$  and  $\mathcal{G}_\pi = \mathcal{G}_{\pi'}$ . Hence,  $\langle \pi', \tau' \rangle$  is our desired DAG-changing walk relative to  $(j, k)$  in  $\mathbf{A}_\mathbf{v}$ .  $\square$

**Lemma C.4** *Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , consider any pair  $\pi^1, \tau^1 \in \Pi(\mathbf{v})$  such that  $\pi^1 = \langle \delta_1, j, k, \delta_2 \rangle$  for some sub-sequences  $\delta_1, \delta_2$  of  $\pi^1$ , and  $\tau^1 = \langle \zeta_1, j, k, \zeta_2 \rangle$  for some sub-sequences  $\zeta_1, \zeta_2$  of  $\tau^1$ . Further consider  $\pi^2 = \langle \delta_1, k, j, \delta_2 \rangle$  and  $\tau^2 = \langle \zeta_1, k, j, \zeta_2 \rangle$ . If  $\mathcal{G}_{\pi^1} = \mathcal{G}_{\tau^1}$ , then  $\mathcal{G}_{\pi^2} = \mathcal{G}_{\tau^2}$ .*

*Proof.* Notice that  $\mathcal{G}_{\pi^2} \in \text{SGS}(\mathcal{P})$  by **Theorem 3.5**. If we can show that  $\tau^2$  is a causal order of  $\mathcal{G}_{\pi^2}$ , it follows from **Lemma B.2** that  $\mathcal{G}_{\pi^2} = \mathcal{G}_{\tau^2}$ . To do so, it suffices to show the following. For any  $i \in \mathbf{v} \setminus \{j, k\}$ ,

- (i) if  $(i \rightarrow j) \in \mathbf{E}(\mathcal{G}_{\pi^2})$ , then  $i \in \zeta_1$ ;
- (ii) if  $(i \rightarrow k) \in \mathbf{E}(\mathcal{G}_{\pi^2})$ , then  $i \in \zeta_1$ ;
- (iii) if  $(j \rightarrow i) \in \mathbf{E}(\mathcal{G}_{\pi^2})$ , then  $i \in \zeta_2$ ;
- (iv) if  $(k \rightarrow i) \in \mathbf{E}(\mathcal{G}_{\pi^2})$ , then  $i \in \zeta_2$ .

For (i), suppose that  $(i \rightarrow j) \in \mathbf{E}(\mathcal{G}_{\pi^2})$ . If  $(i \rightarrow j) \in \mathbf{E}(\mathcal{G}_{\pi^1})$  as well, then  $(i \rightarrow j) \in \mathbf{E}(\mathcal{G}_{\tau^1})$  since  $\mathcal{G}_{\pi^1} = \mathcal{G}_{\tau^1}$ . This entails that  $i \in \zeta_1$ . On the other hand, consider the case that  $(i \rightarrow j) \notin \mathbf{E}(\mathcal{G}_{\pi^1})$ . Then

$$\begin{aligned} X_i \not\perp_{\mathcal{P}} X_j \mid \mathbf{X}_{\delta_1 \setminus \{i\}} \cup \{X_k\} & \quad \because (i \rightarrow j) \in \mathbf{E}(\mathcal{G}_{\pi^2}) & (4) \\ X_i \not\perp_{\mathcal{P}} \{X_j, X_k\} \mid \mathbf{X}_{\delta_1 \setminus \{i\}} & \quad \because (4), \text{weak union} & (5) \\ X_i \perp_{\mathcal{P}} X_j \mid \mathbf{X}_{\delta_1 \setminus \{i\}} & \quad \because (i \rightarrow j) \notin \mathbf{E}(\mathcal{G}_{\pi^1}) & (6) \\ X_i \not\perp_{\mathcal{P}} X_k \mid \mathbf{X}_{\delta_1 \setminus \{i\}} \cup \{X_j\} & \quad \because (5), (6), \text{contraction} & (7) \end{aligned}$$

By (RU), (8) entails that  $(i \rightarrow k) \in \mathbf{E}(\mathcal{G}_{\pi^1}) = \mathbf{E}(\mathcal{G}_{\tau^1})$ . Since  $\tau^1$  is a causal order of  $\mathcal{G}_{\tau^1}$ , we have  $i \in \zeta_1$ .

For (ii), suppose that  $(i \rightarrow k) \in \mathbf{E}(\mathcal{G}_{\pi^2})$ . Similar to (i), the case for  $(i \rightarrow k) \in \mathbf{E}(\mathcal{G}_{\pi^1})$  is simple. So consider the case where  $(i \rightarrow k) \notin \mathbf{E}(\mathcal{G}_{\pi^1})$ .

$$\begin{aligned} X_i \not\perp_{\mathcal{P}} X_k \mid \mathbf{X}_{\delta_1 \setminus \{i\}} & \quad \because (i \rightarrow k) \in \mathbf{E}(\mathcal{G}_{\pi^2}) & (8) \\ X_i \not\perp_{\mathcal{P}} \{X_j, X_k\} \mid \mathbf{X}_{\delta_1 \setminus \{i\}} & \quad \because (8), \text{decomposition} & (9) \\ X_i \perp_{\mathcal{P}} X_k \mid \mathbf{X}_{\delta_1 \setminus \{i\}} \cup \{X_j\} & \quad \because (i \rightarrow k) \notin \mathbf{E}(\mathcal{G}_{\pi^1}) & (10) \\ X_i \not\perp_{\mathcal{P}} X_j \mid \mathbf{X}_{\delta_1 \setminus \{i\}} & \quad \because (9), (10), \text{contraction} & (11) \end{aligned}$$

By (RU), (13) entails that  $(i \rightarrow j) \in \mathbf{E}(\mathcal{G}_{\pi^1}) = \mathbf{E}(\mathcal{G}_{\tau^1})$  and hence  $i \in \zeta_1$ .

For (iii), suppose that  $(j \rightarrow i) \in \mathbf{E}(\mathcal{G}_{\pi^2})$ . Then we have  $(j \rightarrow i) \in \mathbf{E}(\mathcal{G}_{\pi^1})$  by (RU) because  $\text{Pre}(i, \pi^1) = \text{Pre}(i, \pi^2)$ . Hence  $(j \rightarrow i) \in \mathbf{E}(\mathcal{G}_{\tau^1})$  since  $\mathcal{G}_{\pi^1} = \mathcal{G}_{\tau^1}$ . Given that  $\tau^1$  is a causal order of  $\mathcal{G}_{\tau^1}$ , we have  $i \in \zeta_2$ . (iv) is analogous to (iii).  $\square$

**Lemma C.5** *Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , consider any two DAG-changing walks  $\mathfrak{W} = \langle \pi^1, \dots, \pi^m \rangle$  and  $\mathfrak{W}' = \langle \tau^1, \dots, \tau^n \rangle$  in  $\mathbf{A}_\mathbf{v}$  where  $\pi^1 = \tau^1$ . If  $\mathfrak{W}$  and  $\mathfrak{W}'$  are both relative to the same  $(j, k)$  for some  $j, k \in \mathbf{v}$ , then  $\mathcal{G}_{\pi^m} = \mathcal{G}_{\tau^n}$ .*

*Proof.* Immediate from **Definition C.1** and **Lemma C.4**.  $\square$

**Lemma C.6** *Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , consider any DAG-changing walk  $\mathfrak{W} = \langle \pi^1, \dots, \pi^m \rangle$  in  $\mathbf{A}_\mathbf{v}$  which is relative to  $(j, k)$  for some  $j, k \in \mathbf{v}$ . Then  $j \rightarrow k$  is a singular edge in  $\mathcal{G}_{\pi^1}$ .*

*Proof.* Let  $\mathfrak{W}_0$  denotes the DAG-preserving walk  $\langle \pi^1, \dots, \pi^{m-1} \rangle$ . Given that  $\pi^m$  is  $(j, k)$ -different from  $\pi^{m-1}$ , it follows from **Lemma B.3** and (RU) that  $(j \rightarrow k) \in \mathbf{E}(\mathcal{G}_{\pi^{m-1}})$ . Since  $\mathfrak{W}_0$  is a DAG-preserving walk in  $\mathbf{A}_\mathbf{v}$ , we have  $(j \rightarrow k) \in \mathbf{E}(\mathcal{G}_{\pi^1}) = \mathbf{E}(\mathcal{G}_{\pi^{m-1}})$ .

Next, suppose by reductio that  $j \rightarrow k$  is not a singular edge in  $\mathcal{G}_{\pi^1}$ . Then there is a directed path from  $j$  to  $k$  other than  $j \rightarrow k$  in  $\mathcal{G}_{\pi^1}$ . So there exists  $l \in \mathbf{v}$  such that  $l \in \text{De}(j, \mathcal{G}_{\pi^1}) \cap \text{An}(k, \mathcal{G}_{\pi^1})$ . In order to ensure that  $j$  and  $k$  are  $\pi^{m-1}$ -adjacent, either  $\pi^{m-1}[l] < \pi^{m-1}[j]$  or  $\pi^{m-1}[l] > \pi^{m-1}[k]$  holds. However, either case will violate that  $\pi^{m-1}$  is a causal order of  $\mathcal{G}_{\pi^{m-1}} = \mathcal{G}_{\pi^1}$ .  $\square$

**Lemma C.7** Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , consider  $\pi \in \Pi(\mathbf{v})$  where  $(j \rightarrow k) \in \mathbf{E}(\mathcal{G}_\pi)$  is a singular edge for some  $j, k \in \mathbf{v}$ . Then there exists a DAG-changing walk  $\mathfrak{W} = \langle \pi, \dots, \tau \rangle$  in  $\mathbf{A}_\mathbf{v}$  relative to  $(j, k)$  where  $\tau = \text{tuck}(\pi, j, k)$ .

*Proof.* First, we rewrite  $\pi = \langle \delta_{<j}, j, \delta_{j \sim k}, k, \delta_{>k} \rangle$  as usual. Then we partition  $\delta_{j \sim k}$  as follows:  $\zeta_a = \langle i \in \delta_{j \sim k} : i \in \text{An}(k, \mathcal{G}_\pi) \rangle$ , and  $\zeta_b = \langle i \in \delta_{j \sim k} : i \notin \text{An}(k, \mathcal{G}_\pi) \rangle$ . Given that  $(j \rightarrow k)$  is a singular edge, we know that  $\text{De}(j, \mathcal{G}_{\pi^1}) \cap \text{An}(k, \mathcal{G}_{\pi^1}) = \emptyset$ . In other words, we know that (i) each vertex in  $\zeta_a$  has no ancestor in  $\delta_{j \sim k} \setminus \zeta_a$  in  $\mathcal{G}_\pi$  and (ii) each vertex in  $\zeta_b$  has no descendant in  $\delta_{j \sim k} \setminus \zeta_b$  in  $\mathcal{G}_\pi$ .

Now consider the permutation  $\tau' = \langle \delta_{<j}, \zeta_a, j, k, \zeta_b, \delta_{>k} \rangle$  in particular. We want to show that there exists a DAG-preserving walk from  $\pi$  to  $\tau'$ . Such a walk is easy to construct. First, perform repeated ATs by moving each  $i \in \zeta_a$  prior to  $j$  from left to right, and then repeated ATs by moving each  $i \in \zeta_b$  behind  $k$  from right to left. The two sets of ATs are licensed by (i) and (ii) respectively. Hence, we have  $\mathcal{G}_{\tau'} = \mathcal{G}_\pi$ . Finally, consider  $\tau = \text{tuck}(\pi, j, k) = \langle \delta_{<j}, \zeta_a, k, j, \zeta_b, \delta_{>k} \rangle$  which is  $(j, k)$ -different from  $\tau'$ . By (RU) and **Lemma B.3**, we know that  $\mathcal{G}_\tau \neq \mathcal{G}_{\tau'}$  and thus  $\langle \pi, \dots, \tau', \tau \rangle$  is a DAG-changing walk in  $\mathbf{A}_\mathbf{v}$  relative to  $(j, k)$ .  $\square$

**Theorem C.8** Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , consider any DAG-changing walk  $\mathfrak{W} = \langle \pi^1, \dots, \pi^m \rangle$  in  $\mathbf{A}_\mathbf{v}$  which is relative to  $(j, k)$  for some  $j, k \in \mathbf{v}$ . Then  $\mathcal{G}_{\pi^m} = \mathcal{G}_\tau$  where  $\tau = \text{tuck}(\pi^1, j, k)$ .

*Proof.* We obtain a DAG-changing walk  $\mathfrak{W}' = \langle \pi^1, \dots, \tau \rangle$  in  $\mathbf{A}_\mathbf{v}$  relative to  $(j, k)$  by **Lemma C.7**. Since both  $\mathfrak{W}$  and  $\mathfrak{W}'$  are relative to the same  $(j, k)$ , it follows from **Lemma C.5** that  $\mathcal{G}_{\pi^m} = \mathcal{G}_\tau$ .  $\square$

Similar to the discussion in Appendix B, we want to fix the ordering of the set of singular edges in any DAG. This ensures that ESP and unbounded GRaSP<sub>1</sub> will not yield different DAGs simply due to the issue of order-dependence. Below we prove that ESP and unbounded GRaSP<sub>1</sub> are equivalent algorithms.

**Theorem C.9** Given a graphoid  $\mathcal{P}$  and any initial permutation  $\pi \in \Pi(\mathbf{v})$ , the DAG induced by the output of unbounded GRaSP<sub>1</sub> is equivalent to the DAG returned by ESP.

*Proof.* Consider any  $j, k \in \mathbf{v}$ . By **Lemma C.6** and **Lemma C.7**, every DAG-changing walk  $\mathfrak{W} = \langle \pi^1, \dots, \pi^m \rangle$  in  $\mathbf{A}_\mathbf{v}$  relative to  $(j, k)$  corresponds to a *tuck* operation of the singular edge  $j \rightarrow k$  in  $\mathbf{E}(\mathcal{G}_{\pi^1})$ . Hence, by **Lemma C.3**, we know that  $\text{tuck}(\pi^1, j, k)$  corresponds to the neighboring relation between  $\mathcal{G}_{\pi^1}$  and  $\mathcal{G}_{\pi^m}$  in  $\mathbf{A}_\mathbf{v}(\mathcal{P})$  that are  $(j, k)$ -reverse. Therefore, every step taken by ESP to move to a neighboring state in  $\mathbf{A}_\mathbf{v}(\mathcal{P})$  (relative to a unique pair of vertices) is equivalent to the tuck operation taken by GRaSP<sub>1</sub> over the same pair of vertices.  $\square$

## D CAUSAL RAZORS AND GRASP

In this section, we first provide a logical analysis of the causal razors discussed in the main text.<sup>5</sup> Then we construct new causal razors with respect to each tier of GRaSP, and show how a higher tier of GRaSP requires a strictly weaker causal razor.

**Theorem D.1** The following statements are true:

- (a) For any joint probability distribution  $\mathcal{P}$ ,  $\text{uPm}(\mathcal{P}) = \text{CFC}(\mathcal{P}) \subseteq \text{uFr}(\mathcal{P}) \subseteq \text{Fr}(\mathcal{P}) \subseteq \text{Pm}(\mathcal{P}) \subseteq \text{SGS}(\mathcal{P})$ .
- (b) For any joint probability distribution  $\mathcal{P}$ , if faithfulness is satisfied,  $\text{CFC}(\mathcal{P}) = \text{uFr}(\mathcal{P}) = \text{Fr}(\mathcal{P}) = \text{Pm}(\mathcal{P})$ .
- (c) There exists a joint probability distribution s.t.  $\text{CFC}(\mathcal{P}) \subset \text{uFr}(\mathcal{P})$ .
- (d) There exists a joint probability distribution s.t.  $\text{uFr}(\mathcal{P}) \subset \text{Fr}(\mathcal{P})$ .
- (e) There exists a joint probability distribution s.t.  $\text{Fr}(\mathcal{P}) \subset \text{Pm}(\mathcal{P})$ .
- (f) There exists a joint probability distribution s.t.  $\text{Pm}(\mathcal{P}) \subset \text{SGS}(\mathcal{P})$ .

<sup>5</sup>There are other causal razors discussed in the literature, including, but not limited to, *adjacency-faithfulness* and *orientation-faithfulness* in [Ramsey et al., 2006], and *triangle-faithfulness* in Zhang [2013]. But they do not have a strong connection with our discussion of GRaSP and so will not be analyzed in this work.

*Proof.* For (a),  $\text{uPm}(\mathcal{P}) = \text{CFC}(\mathcal{P})$  is our result in **Theorem B.15**.  $\text{CFC}(\mathcal{P}) \subseteq \text{uFr}(\mathcal{P})$  is proven in [Raskutti and Uhler, 2018],  $\text{uFr}(\mathcal{P}) \subseteq \text{Fr}(\mathcal{P})$  is true by **Definition 3.3**,  $\text{Fr}(\mathcal{P}) \subseteq \text{Pm}(\mathcal{P})$  in [Forster et al., 2020], and  $\text{Pm}(\mathcal{P}) \subseteq \text{SGS}(\mathcal{P})$  in [Zhang, 2013]. (b) is a direct consequence of (a) and **Theorem 3.8**.

For (c), see [Raskutti and Uhler, 2018, Theorem 2.4]. For (d), see [Forster et al., 2020, Figure 6]. For (e), see [Raskutti and Uhler, 2018, Theorem 2.5]. For (f), see [Zhang, 2013, Figure 2]. Additionally, the example in **Theorem D.6** and its corresponding Figure 2 verifies (c) and (e):  $\mathcal{G}^* \in \text{uFr}(\mathcal{P}) \setminus \text{CFC}(\mathcal{P})$  and  $\mathcal{G}_\pi \in \text{Pm}(\mathcal{P}) \setminus \text{Fr}(\mathcal{P})$ . On the other hand, each of  $\mathcal{G}_{\pi^1}, \mathcal{G}_{\pi^2}, \mathcal{G}_{\pi^3}$ , and  $\mathcal{G}_{\pi^4}$  in the DAG-associahedron in Figure 1 is in  $\text{SGS}(\mathcal{P}) \setminus \text{Pm}(\mathcal{P})$  verifying (f).  $\square$

**Definition D.2** (*TSP-razor and ESP-razor*) Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , let  $\text{tsp}(\mathcal{P}, \pi)$  be the DAG returned by TSP on  $\mathcal{P}$  by setting  $\pi$  as the initial permutation. Define  $\text{TSP}(\mathcal{P}) = \{\mathcal{G} \in \text{DAG}(\mathbf{V}) : \pi \in \Pi(\mathbf{v}) \text{ and } \mathcal{G} = \text{tsp}(\mathcal{P}, \pi)\}$  as the set of DAGs returned by TSP on  $\mathcal{P}$  over each initial permutation in  $\Pi(\mathbf{v})$ . Further define

$$\text{TSPr}(\mathcal{P}) = \{\mathcal{G} \in \text{TSP}(\mathcal{P}) : \neg \exists \mathcal{G}' \in \text{TSP}(\mathcal{P}) \text{ s.t. } \mathcal{G}' \notin \text{MEC}(\mathcal{G})\}.$$

$(\mathcal{G}^*, \mathcal{P})$  satisfies the TSP-razor if  $\mathcal{G}^* \in \text{TSPr}(\mathcal{P})$ . Similarly for ESP, esp, ESP, ESPr, and ESP-razor.

One can observe that  $\text{TSPr}(\mathcal{P}) = \text{TSP}(\mathcal{P})$  if every DAG in  $\text{TSP}(\mathcal{P})$  belongs to the same MEC, and  $\text{TSPr}(\mathcal{P}) = \emptyset$  otherwise. The same is also true for  $\text{ESPr}(\mathcal{P})$  and  $\text{ESP}(\mathcal{P})$ . These definitions will be proven useful when we compare them with the classes of DAGs discussed in **Theorem D.1**. Below we provide a similar definition for each tier of GRaSP.

**Definition D.3** (*GRaSP<sub>t</sub>-razor*) Given a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ , for  $t \in \{0, 1, 2\}$ , define  $\text{GRaSP}_t(\mathcal{P}) = \{\mathcal{G}_\tau \in \text{DAG}(\mathbf{V}) : \pi \in \Pi(\mathbf{v}) \text{ and } \tau = \text{grasp}(\mathcal{P}, \pi, |\mathbf{v}|, t)\}$  as the set of DAGs returned by unbounded GRaSP<sub>t</sub> on  $\mathcal{P}$  over each initial permutation in  $\Pi(\mathbf{v})$ . Further define

$$\text{GRaSP}_{t\text{r}}(\mathcal{P}) = \{\mathcal{G} \in \text{GRaSP}_t(\mathcal{P}) : \neg \exists \mathcal{G}' \in \text{GRaSP}_t(\mathcal{P}) \text{ s.t. } \mathcal{G}' \notin \text{MEC}(\mathcal{G})\}.$$

$(\mathcal{G}^*, \mathcal{P})$  satisfies the GRaSP<sub>t</sub>-razor if  $\mathcal{G}^* \in \text{GRaSP}_{t\text{r}}(\mathcal{P})$ .

**Theorem D.4** Given a graphoid  $\mathcal{P}$ , the following statement is true:

$$\text{CFC}(\mathcal{P}) = \text{TSPr}(\mathcal{P}) = \text{GRaSP}_{0\text{r}}(\mathcal{P}) \subseteq \text{ESPr}(\mathcal{P}) = \text{GRaSP}_{1\text{r}}(\mathcal{P}) \subseteq \text{GRaSP}_{2\text{r}}(\mathcal{P}) \subseteq \text{uFr}(\mathcal{P}).$$

*Proof.*  $\text{CFC}(\mathcal{P}) = \text{TSPr}(\mathcal{P}) = \text{GRaSP}_{0\text{r}}(\mathcal{P})$  is directly entailed by **Theorem 4.7** and **Theorem 4.8**. Solus et al. [2021] showed that  $\text{TSPr}(\mathcal{P}) \subseteq \text{ESPr}(\mathcal{P})$ .  $\text{ESPr}(\mathcal{P}) = \text{GRaSP}_{1\text{r}}(\mathcal{P})$  is entailed by **Theorem C.9**.

Next, to show that  $\text{GRaSP}_{1\text{r}}(\mathcal{P}) \subseteq \text{GRaSP}_{2\text{r}}(\mathcal{P})$ , notice that  $\text{GRaSP}_{1\text{r}}(\mathcal{P}) = \text{GRaSP}_1(\mathcal{P})$  when all DAGs in  $\text{GRaSP}_1(\mathcal{P})$  belong to the same MEC, and  $\text{GRaSP}_{1\text{r}}(\mathcal{P}) = \emptyset$  otherwise. The latter case validates  $\text{GRaSP}_{1\text{r}}(\mathcal{P}) \subseteq \text{GRaSP}_{2\text{r}}(\mathcal{P})$  trivially. Now consider the former case where all DAGs in the non-empty  $\text{GRaSP}_1(\mathcal{P})$  belong to the same MEC and so they have the same number of edges. Now consider any  $\pi \in \Pi(\mathbf{v})$  satisfying  $\mathcal{G}_\pi \in \text{Fr}(\mathcal{P})$  (where  $\text{Fr}(\mathcal{P})$  is necessarily non-empty). We know that  $\mathcal{G}_\pi \in \text{GRaSP}_1(\mathcal{P})$ . This is because every initial permutation in  $\Pi(\mathbf{v})$  is considered and unbounded GRaSP<sub>1</sub> will never return a denser permutation than its initial permutation. Hence, every DAG in  $\text{GRaSP}_1(\mathcal{P})$  is the sparsest Markovian DAG. (The same also holds when  $\text{GRaSP}_{2\text{r}}(\mathcal{P}) \neq \emptyset$ .) Then the construction of **Algorithm 2** entails that GRaSP<sub>2</sub> will return the same permutation as GRaSP<sub>1</sub>. Hence,  $\text{GRaSP}_{1\text{r}}(\mathcal{P}) = \text{GRaSP}_{2\text{r}}(\mathcal{P})$  when all DAGs in  $\text{GRaSP}_1(\mathcal{P})$  belongs to the same MEC.

Lastly, to show that  $\text{GRaSP}_{2\text{r}}(\mathcal{P}) \subseteq \text{uFr}(\mathcal{P})$ , we use a proof similar to the above. First, the case where  $\text{GRaSP}_{2\text{r}}(\mathcal{P}) = \emptyset$  is trivial. Consider the case where  $\text{GRaSP}_{2\text{r}}(\mathcal{P}) = \text{GRaSP}_2(\mathcal{P})$  s.t. all DAGs in  $\text{GRaSP}_2(\mathcal{P})$  are in the same MEC. Using a similar inference used in the last paragraph, we know that every DAG in  $\text{GRaSP}_2(\mathcal{P})$  is the sparsest Markovian DAG. Therefore,  $\text{GRaSP}_{2\text{r}}(\mathcal{P}) = \text{uFr}(\mathcal{P})$  when all DAGs in  $\text{GRaSP}_2(\mathcal{P})$  belongs to the same MEC.  $\square$

**Theorem D.5** There exists a graphoid  $\mathcal{P}$  s.t.  $\text{GRaSP}_{0\text{r}}(\mathcal{P}) \subset \text{GRaSP}_{1\text{r}}(\mathcal{P})$ .

*Proof.* Given the equivalence between TSP and unbounded GRaSP<sub>0</sub> shown in **Theorem 4.7**, and that between ESP and unbounded GRaSP<sub>1</sub> in **Theorem C.9**, we can borrow the example from [Solus et al., 2021] on how ESP requires a strictly weaker causal razor than TSP. We refer the readers to Figure 3 in the [supplementary materials](#) of [Solus et al., 2021].  $\square$

In the remainder of this section, we discuss two examples: how unbounded  $\text{GRaSP}_2$  requires a strictly weaker causal razor than unbounded  $\text{GRaSP}_1$ , and how unbounded  $\text{GRaSP}_2$  requires a strictly stronger causal razor than u-frugality. The joint distribution of each example below is a compositional graphoid. For the sake of simplicity, we only include CI relations that hold between two *singleton* sets of variables such that all other CI relations entailed by each of the graphoid axioms discussed in Appendix A.2 are understood.

**Theorem D.6** *There exists a graphoid  $\mathcal{P}$  s.t.  $\text{GRaSP}_{1r}(\mathcal{P}) \subset \text{GRaSP}_{2r}(\mathcal{P})$ .*

*Proof.* Given  $\mathbf{V} = \{X_1, \dots, X_4\}$ , consider the unfaithful model  $(\mathcal{G}^*, \mathcal{P})$  where the true DAG  $\mathcal{G}^*$  is shown on the left in Figure 2, and  $\mathbf{I}(\mathcal{P}) = \Phi \cup \Psi$  where  $\Phi$  is the set of faithful CI relations and  $\Psi$  is the set of unfaithful CI relations as listed below:

$$\begin{aligned} \Phi &= \{\phi_1 = \langle X_1, X_3 \mid \{X_2\} \rangle, \phi_2 = \langle X_2, X_4 \mid \{X_1, X_3\} \rangle\}; \\ \Psi &= \{\psi_1 = \langle X_2, X_4 \mid \emptyset \rangle\}. \end{aligned}$$

For every  $\mathcal{G} \in \text{CMC}(\mathcal{P})$  where  $\psi_1 \in \mathbf{I}(\mathcal{G})$ , we have  $5 = |\mathbf{E}(\mathcal{G})| > |\mathbf{E}(\mathcal{G}^*)| = 4$ . Also, all 4-edge Markovian DAGs are in the same MEC. Hence, u-frugality is satisfied. Consider feeding the initial permutation  $\pi = \langle 2, 4, 1, 3 \rangle$  to unbounded  $\text{GRaSP}_1$ . It will return the same  $\pi$  after the DFS procedure and the induced  $\mathcal{G}_\pi$ , as shown on the right in Figure 2, contains 5 edges. Therefore, unbounded  $\text{GRaSP}_1$  fails to return the sparsest permutation under some initial permutation and  $\text{GRaSP}_{1r}(\mathcal{P}) = \emptyset$ .

On the contrary,  $|\mathbf{v}|! = 24$  initial permutations have been tested on unbounded  $\text{GRaSP}_2$  and it returns  $\hat{\tau}$  where  $\mathcal{G}_{\hat{\tau}} \in \text{MEC}(\mathcal{G}^*)$  for each initial permutation. Hence,  $\text{GRaSP}_{2r}(\mathcal{P}) \neq \emptyset$ .  $\square$



Figure 2: An unfaithful model satisfying u-frugality. The true DAG  $\mathcal{G}^*$  is shown on the left where the two shaded vertices indicate the unfaithful marginal independence  $X_2 \perp\!\!\!\perp_{\mathcal{P}} X_4 \mid \emptyset$ . Unbounded  $\text{GRaSP}_1$  returns its initial permutation  $\pi = \langle 2, 4, 1, 3 \rangle$ . The induced DAG  $\mathcal{G}_\pi$  is shown on the right with 5 edges. However, unbounded  $\text{GRaSP}_2$  manages to return one of the sparsest permutations under every initial permutation.

**Theorem D.7** *There exists a graphoid  $\mathcal{P}$  s.t.  $\text{GRaSP}_{2r}(\mathcal{P}) \subset \text{uFr}(\mathcal{P})$ .*

*Proof.* The example below is one of the uDAGs studied in Table 1 in Section 5.1 where  $\text{GRaSP}_2$  fails to return one of the sparsest permutations under u-frugality. Given  $\mathbf{V} = \{X_1, \dots, X_5\}$ , consider the unfaithful model  $(\mathcal{G}^*, \mathcal{P})$  where the true DAG  $\mathcal{G}^*$  is shown on the left in Figure 3, and  $\mathbf{I}(\mathcal{P}) = \Phi \cup \Psi$  where  $\Phi$  is the set of faithful CI relations and  $\Psi$  is the set of unfaithful CI relations as listed below:

$$\begin{aligned} \Phi &= \{\phi_1 = \langle X_1, X_2 \mid \emptyset \rangle, \phi_2 = \langle X_1, X_2 \mid \{X_3\} \rangle, \\ &\quad \phi_3 = \langle X_2, X_3 \mid \emptyset \rangle, \phi_4 = \langle X_2, X_3 \mid \{X_1\} \rangle, \\ &\quad \phi_5 = \langle X_2, X_5 \mid \{X_1, X_3, X_4\} \rangle\}; \\ \Psi &= \{\psi_1 = \langle X_1, X_5 \mid \emptyset \rangle\}. \end{aligned}$$

For every  $\mathcal{G} \in \text{CMC}(\mathcal{P})$  where  $\psi_1 \in \mathbf{I}(\mathcal{G})$ , we have  $|\mathbf{E}(\mathcal{G})| > |\mathbf{E}(\mathcal{G}^*)| = 7$ . Also, all 7-edge Markovian DAGs are in the same MEC and there exists no sparser Markovian DAG. Hence, u-frugality is satisfied s.t.  $\text{uFr}(\mathcal{P}) \neq \emptyset$ .

Next, consider feeding the initial permutation  $\pi = \langle 5, 1, 3, 4, 2 \rangle$  to unbounded  $\text{GRaSP}_2$ . It will return the same  $\pi$  after the DFS procedure and the induced  $\mathcal{G}_\pi$ , as shown on the right in Figure 3, contains 8 edges. Therefore, unbounded  $\text{GRaSP}_2$  fails to return one of the sparsest permutations under some initial permutation and  $\text{GRaSP}_{2r}(\mathcal{P}) = \emptyset$ .  $\square$



Figure 3: An unfaithful model satisfying u-frugality. The true DAG  $\mathcal{G}^*$  is shown on the left where the two shaded vertices indicate the unfaithful marginal independence  $X_1 \perp\!\!\!\perp X_5 \mid \emptyset$ . Unbounded GRaSP<sub>2</sub> returns its initial permutation  $\pi = \langle 5, 1, 3, 4, 2 \rangle$ . The induced DAG  $\mathcal{G}_\pi$  is shown on the right with 8 edges. Hence, GRaSP<sub>2</sub> is not correct under u-frugality alone.

**Corollary 4.9** *Given a graphoid  $\mathcal{P}$ , unbounded GRaSP<sub>2</sub> is correct under a strictly weaker causal razor than unbounded GRaSP<sub>1</sub>, which is correct under a strictly weaker causal razor than unbounded GRaSP<sub>0</sub>.*

## E GROW-SHRINK ALGORITHM AND ITS PROPERTIES

**Definition E.1** *Given an observational dataset  $\mathcal{D}$  with  $n$  i.i.d. observations from a joint probability distribution  $\mathcal{P}$  over  $\mathbf{V}$  that belongs to a curved exponential family<sup>6</sup>, for every  $X \in \mathbf{V}$  and every  $\mathbf{M} \subseteq \mathbf{V} \setminus X$ ,*

$$\text{BIC}_{\mathcal{D}}(X, \mathbf{M}) = \ell_{X|\mathbf{M}}(\hat{\theta}_{\text{MLE}} \mid \mathcal{D}) + c \frac{|\hat{\theta}_{\text{MLE}}|}{2} \log(n)$$

where  $\ell_{X|\mathbf{M}}$  is the conditional log likelihood function,  $|\hat{\theta}_{\text{MLE}}|$  is the absolute value of the maximum likelihood estimate, and  $c$  is a multiplier for the parameter penalty.

BIC score is a *decomposable* scoring function in the sense that the BIC score of any DAG  $\mathcal{G}$  (over the same set of variables  $\mathbf{V}$  as the observational dataset  $\mathcal{D}$ ), denoted as  $\text{BIC}_{\mathcal{D}}(\mathcal{G})$ , satisfies the following:

$$\text{BIC}_{\mathcal{D}}(\mathcal{G}) = \sum_{i \in \mathbf{V}} \text{BIC}_{\mathcal{D}}(X_i, \mathbf{X}_{\text{Pa}(i, \mathcal{G})}).$$

In addition, since we will be using BIC throughout this appendix, we assume that every joint probability distribution  $\mathcal{P}$  belongs to a curved exponential family in this section.

---

**Algorithm 1:** GROW:  $grow(\mathcal{D}, X, \mathbf{Z})$

---

**Input:** (a)  $\mathcal{D}$ : an observational dataset over  $\mathbf{V}$ ; (b)  $X \in \mathbf{V}$ ; (c)  $\mathbf{Z} \subseteq \mathbf{V} \setminus \{X\}$

**Output:**  $\mathbf{M}_{gr} \subseteq \mathbf{Z}$

```

1  $s \leftarrow \text{BIC}_{\mathcal{D}}(X, \emptyset)$ 
2  $s' \leftarrow s$ 
3  $\mathbf{M}_{gr} \leftarrow \emptyset$ 
4 do
5    $s \leftarrow s'$ 
6    $s' \leftarrow \max_{Y \in \mathbf{Z} \setminus \mathbf{M}_{gr}} \text{BIC}_{\mathcal{D}}(X, \mathbf{M}_{gr} \cup \{Y\})$ 
7    $Y' \leftarrow \text{argmax}_{Y \in \mathbf{Z} \setminus \mathbf{M}_{gr}} \text{BIC}_{\mathcal{D}}(X, \mathbf{M}_{gr} \cup \{Y\})$ 
8   if  $s' > s$  then
9      $\mathbf{M}_{gr} \leftarrow \mathbf{M}_{gr} \cup \{Y'\}$ 
10 while  $s' > s$ 
11 return  $\mathbf{M}_{gr}$ 

```

---

<sup>6</sup>See [Kass and Vos, 2011] for an in-depth analysis of curved exponential families.

---

**Algorithm 2:** SHRINK:  $shrink(\mathcal{D}, X, \mathbf{Z})$ 

---

**Input:** (a)  $\mathcal{D}$ : an observational dataset over  $\mathbf{V}$ ; (b)  $X \in \mathbf{V}$ ; (c)  $\mathbf{Z} \subseteq \mathbf{V} \setminus \{X\}$

**Output:** (i)  $\mathbf{M}_{sh} \subseteq \mathbf{Z}$ ; (ii)  $s = \text{BIC}_{\mathcal{D}}(X, \mathbf{M}_{sh})$

```
1  $s \leftarrow \text{BIC}_{\mathcal{D}}(X, \mathbf{Z})$ 
2  $s' \leftarrow s$ 
3  $\mathbf{M}_{sh} \leftarrow \mathbf{Z}$ 
4 do
5    $s \leftarrow s'$ 
6    $s' \leftarrow \max_{Y \in \mathbf{M}_{sh}} \text{BIC}_{\mathcal{D}}(X, \mathbf{M}_{sh} \setminus \{Y\})$ 
7    $Y' \leftarrow \text{argmax}_{Y \in \mathbf{M}_{sh}} \text{BIC}_{\mathcal{D}}(X, \mathbf{M}_{sh} \setminus \{Y\})$ 
8   if  $s' > s$  then
9      $\mathbf{M}_{sh} \leftarrow \mathbf{M}_{sh} \setminus \{Y'\}$ 
10 while  $s' > s$ 
11 return  $\mathbf{M}_{sh}, s$ 
```

---

**Theorem E.2** [Chickering, 2002] Given an observational dataset  $\mathcal{D}$  with  $n$  i.i.d. observations from a joint probability distribution  $\mathcal{P}$  over  $\mathbf{V}$ , consider  $\mathcal{G}, \mathcal{G}' \in \text{DAG}(\mathbf{V})$  where  $\mathcal{G}'$  is resulted from adding the edge  $j \rightarrow k$  in  $\mathcal{G}$ . In the large sample limit of  $n$ ,

- (a) if  $X_j \not\perp_{\mathcal{P}} X_k \mid \mathbf{X}_{\text{Pa}(k, \mathcal{G})}$ , then  $\text{BIC}_{\mathcal{D}}(\mathcal{G}') > \text{BIC}_{\mathcal{D}}(\mathcal{G})$ ;
- (b) if  $X_j \perp_{\mathcal{P}} X_k \mid \mathbf{X}_{\text{Pa}(k, \mathcal{G})}$ , then  $\text{BIC}_{\mathcal{D}}(\mathcal{G}') < \text{BIC}_{\mathcal{D}}(\mathcal{G})$ .

The theorem above is known as the *local consistency* of BIC score over DAGs. We can easily derive a lemma which concerns the BIC score of a variable (relative to a set of variables).

**Lemma E.3** Given an observational dataset  $\mathcal{D}$  with  $n$  i.i.d. observations from a joint probability distribution  $\mathcal{P}$  over  $\mathbf{V}$ , consider any distinct  $j, k \in \mathbf{v}$  and  $\mathbf{i} \subseteq \mathbf{v} \setminus \{j, k\}$ . In the large sample limit of  $n$ ,

- (a) if  $X_j \not\perp_{\mathcal{P}} X_k \mid \mathbf{X}_{\mathbf{i}}$ , then  $\text{BIC}_{\mathcal{D}}(X_k, \mathbf{X}_{\mathbf{i}} \cup \{X_j\}) > \text{BIC}_{\mathcal{D}}(X_k, \mathbf{X}_{\mathbf{i}})$ ;
- (b) if  $X_j \perp_{\mathcal{P}} X_k \mid \mathbf{X}_{\mathbf{i}}$ , then  $\text{BIC}_{\mathcal{D}}(X_k, \mathbf{X}_{\mathbf{i}} \cup \{X_j\}) < \text{BIC}_{\mathcal{D}}(X_k, \mathbf{X}_{\mathbf{i}})$ .

*Proof.* Construct a DAG  $\mathcal{G} \in \text{DAG}(\mathbf{V})$  by drawing all and only directed edges from each vertex in  $\mathbf{i}$  to  $k$ , and another DAG  $\mathcal{G}' \in \text{DAG}(\mathbf{V})$  by adding  $j \rightarrow k$  in  $\mathcal{G}$ . Then the lemma immediately follows from **Theorem E.2** and the decomposable feature of BIC scores.  $\square$

**Lemma E.4** Consider an observational dataset  $\mathcal{D}$  with  $n$  i.i.d. observations from a compositional graphoid  $\mathcal{P}$  over  $\mathbf{V}$ . In the large sample limit of  $n$ , for any  $X \in \mathbf{V}$  and any  $\mathbf{Z} \subseteq \mathbf{V} \setminus \{X\}$ ,  $\text{MB}(X, \mathbf{Z}) \subseteq \mathbf{M}_{gr}$  where  $\mathbf{M}_{gr} = \text{grow}(\mathcal{D}, X, \mathbf{Z}) \subseteq \mathbf{Z}$ .

*Proof.* First, **Algorithm 1** requires that  $\mathbf{M}_{gr} \subseteq \mathbf{Z}$ , and  $\text{BIC}_{\mathcal{D}}(X, \mathbf{M}_{gr} \cup \{Y\}) < \text{BIC}_{\mathcal{D}}(X, \mathbf{M}_{gr})$  for every  $Y \in \mathbf{Z} \setminus \mathbf{M}_{gr}$ . By **Lemma E.3**, we have  $X \perp_{\mathcal{P}} Y \mid \mathbf{M}_{gr}$  for each  $Y \in \mathbf{Z} \setminus \mathbf{M}_{gr}$ . By composition, we have  $X \perp_{\mathcal{P}} (\mathbf{Z} \setminus \mathbf{M}_{gr}) \mid \mathbf{M}_{gr}$ . Therefore, by **Definition A.1** and **Lemma A.2**, we have  $\text{MB}(X, \mathbf{Z}) \subseteq \mathbf{M}_{gr}$ .  $\square$

**Lemma E.5** Consider an observational dataset  $\mathcal{D}$  with  $n$  i.i.d. observations from a graphoid  $\mathcal{P}$  over  $\mathbf{V}$ . In the large sample limit of  $n$ , for any  $X \in \mathbf{V}$  and any  $\mathbf{Z} \subseteq \mathbf{V} \setminus \{X\}$ ,  $\text{MB}(X, \mathbf{Z}) = \mathbf{M}_{sh}$  where  $\mathbf{M}_{sh} = \text{shrink}(\mathcal{D}, X, \mathbf{Z}) \subseteq \mathbf{Z}$ .

*Proof.* We show the lemma by  $\mathbf{M}_{sh} \subseteq \text{MB}(X, \mathbf{Z})$  and  $\mathbf{M}_{sh} \supseteq \text{MB}(X, \mathbf{Z})$ .

$[\subseteq]$  By reductio, suppose that there exists  $Y \in \mathbf{M}_{sh} \subseteq \mathbf{Z}$  but  $Y \notin \text{MB}(X, \mathbf{Z})$ . Let  $\mathbf{S}$  be  $\mathbf{M}_{sh} \setminus \{Y\}$ . **Algorithm 2** requires that  $\text{BIC}_{\mathcal{D}}(X, \mathbf{M}_{sh} \setminus \{Y\}) < \text{BIC}_{\mathcal{D}}(X, \mathbf{M}_{sh})$ . In other words, we have  $\text{BIC}_{\mathcal{D}}(X, \mathbf{S}) < \text{BIC}_{\mathcal{D}}(X, \mathbf{S} \cup \{Y\})$ . By **Lemma E.3**, we have  $X \not\perp_{\mathcal{P}} Y \mid \mathbf{S}$ .

Let  $\mathbf{W} = \mathbf{S} \setminus \text{MB}(X, \mathbf{Z})$ . From  $Y \notin \text{MB}(X, \mathbf{Z})$  and  $Y \notin \mathbf{S}$ , we have  $\{Y\} \cup \mathbf{W} \subseteq \mathbf{Z} \setminus \text{MB}(X, \mathbf{Z})$ . Recall **Definition A.1** that  $X \perp\!\!\!\perp_{\mathcal{P}} \mathbf{Z} \setminus \text{MB}(X, \mathbf{Z}) \mid \text{MB}(X, \mathbf{Z})$ . Thus,

$$X \perp\!\!\!\perp_{\mathcal{P}} \{Y\} \cup \mathbf{W} \mid \text{MB}(X, \mathbf{Z}) \quad \because X \perp\!\!\!\perp_{\mathcal{P}} \mathbf{Z} \setminus \text{MB}(X, \mathbf{Z}) \mid \text{MB}(X, \mathbf{Z}), \text{ decomposition} \quad (12)$$

$$X \perp\!\!\!\perp_{\mathcal{P}} Y \mid \text{MB}(X, \mathbf{Z}) \cup \mathbf{W} \quad \because (12), \text{ weak union} \quad (13)$$

$$X \perp\!\!\!\perp_{\mathcal{P}} Y \mid \mathbf{S} \quad \because (13), \mathbf{W} = \mathbf{S} \setminus \text{MB}(X, \mathbf{Z}) \quad (14)$$

Contradiction arises with  $X \not\perp\!\!\!\perp_{\mathcal{P}} Y \mid \mathbf{S}$ .

[ $\supseteq$ ] Observe that **Algorithm 2** removes one variable in  $\mathbf{Z}$  one at a time repeatedly to form  $\mathbf{M}_{sh}$ . Thus, the shrink-procedure corresponds to a sequence of sets of variables  $\langle \mathbf{M}^0, \dots, \mathbf{M}^k \rangle$  and a sequence of variables  $\mathbf{W} = \langle W_1, \dots, W_k \rangle = \mathbf{Z} \setminus \mathbf{M}_{sh}$  such that  $\mathbf{M}^0 = \mathbf{Z}$ ,  $\mathbf{M}^k = \mathbf{M}_{sh}$ , and  $\mathbf{M}^i = \mathbf{M}^{i-1} \setminus \{W_i\}$  (where  $W_i \in \mathbf{M}^{i-1}$ ) for each  $1 < i \leq k$ .

Notice that  $\mathbf{M}^{i-1} = \mathbf{M}^i \cup \{W_i\}$ . **Algorithm 2** requires that  $\text{BIC}_{\mathcal{D}}(X, \mathbf{M}^i) > \text{BIC}_{\mathcal{D}}(X, \mathbf{M}^{i-1}) = \text{BIC}_{\mathcal{D}}(X, \mathbf{M}^i \cup \{W_i\})$ . We then have

$$X \perp\!\!\!\perp_{\mathcal{P}} W_1 \mid \mathbf{M}^1 \quad \because \text{BIC}_{\mathcal{D}}(X, \mathbf{M}^1) > \text{BIC}_{\mathcal{D}}(X, \mathbf{M}^0), \text{ Lemma E.3} \quad (15)$$

$$X \perp\!\!\!\perp_{\mathcal{P}} W_2 \mid \mathbf{M}^2 \quad \because \text{BIC}_{\mathcal{D}}(X, \mathbf{M}^2) > \text{BIC}_{\mathcal{D}}(X, \mathbf{M}^1), \text{ Lemma E.3} \quad (16)$$

$$X \perp\!\!\!\perp_{\mathcal{P}} W_1 \mid \mathbf{M}^2 \cup \{W_2\} \quad \because (15), \mathbf{M}^1 = \mathbf{M}^2 \cup \{W_2\} \quad (17)$$

$$X \perp\!\!\!\perp_{\mathcal{P}} \{W_1, W_2\} \mid \mathbf{M}^2 \quad \because (16), (17), \text{ contraction} \quad (18)$$

$\vdots$

$$X \perp\!\!\!\perp_{\mathcal{P}} \{W_1, \dots, W_k\} \mid \mathbf{M}^k \quad \because \dots, \text{ contraction} \quad (19)$$

$$X \perp\!\!\!\perp_{\mathcal{P}} \mathbf{W} \mid \mathbf{M}_{sh} \quad \because (19), \mathbf{W} = \langle W_1, \dots, W_k \rangle \text{ and } \mathbf{M}^k = \mathbf{M}_{sh} \quad (20)$$

$$X \perp\!\!\!\perp_{\mathcal{P}} \mathbf{Z} \setminus \mathbf{M}_{sh} \mid \mathbf{M}_{sh} \quad \because (20), \mathbf{W} = \mathbf{Z} \setminus \mathbf{M}_{sh} \quad (21)$$

Hence, it follows from **Definition A.1** that  $\mathbf{M}_{sh} \supseteq \text{MB}(X, \mathbf{Z})$ .  $\square$

**Theorem E.6** Consider an observational dataset  $\mathcal{D}$  with  $n$  i.i.d. observations from a compositional graphoid  $\mathcal{P}$  over  $\mathbf{V}$ . In the large sample limit of  $n$ , for any  $X \in \mathbf{V}$  and any  $\mathbf{Z} \subseteq \mathbf{V} \setminus \{X\}$ ,  $\text{MB}(X, \mathbf{Z}) = \mathbf{M}_{gs}$  where  $\mathbf{M}_{gs} = \text{shrink}(\mathcal{D}, X, \text{grow}(\mathcal{D}, X, \mathbf{Z}))$ .

*Proof.* Immediate from **Lemma E.4** and **Lemma E.5**.  $\square$

**Theorem E.7** Consider an observational dataset  $\mathcal{D}$  with  $n$  i.i.d. observations from a (compositional) graphoid  $\mathcal{P}$  over  $\mathbf{V} = \{X_1, \dots, X_m\}$ , and any  $\pi \in \Pi(\mathbf{v})$ . Let  $s_i$  and  $\mathbf{M}_i$  be the score and the set of variables returned by  $\text{shrink}(\mathcal{D}, X_i, \mathbf{X}_{\text{Pre}(i, \pi)})$  (or  $\text{shrink}(\mathcal{D}, X_i, \text{grow}(\mathcal{D}, X_i, \mathbf{X}_{\text{Pre}(i, \pi)}))$ ) if  $\mathcal{P}$  is a compositional graphoid) respectively. Denote  $s_{\pi}$  as  $\sum_{i \in \mathbf{v}} s_i$ . In the large sample limit of  $n$ ,  $\text{BIC}_{\mathcal{D}}(\mathcal{G}_{\pi}) = s_{\pi}$  where  $\mathcal{G}_{\pi}$  is induced from  $\pi$  by (VP).

*Proof.* Immediate from the decomposable feature of BIC scores, **Lemma E.5** and **Theorem E.6**.  $\square$

Lastly, though a compositional graphoid is a sufficient condition for the correct identification of the unique Markov boundary using the grow-shrink algorithm, we are aware of an assumption weaker than compositional graphoid to validate such an identification. Nevertheless, this discussion will be beyond the scope of this paper and we will leave the formal proof to future work.

## F ADDITIONAL EXAMPLES

### F.1 LU ET AL. COMPARISON

Reported below are average statistics obtained by running GRASP<sub>2</sub> on the published datasets used to generate Figure 6 in [Lu et al., 2021]<sup>7</sup>. We cannot compare these results to Lu et al. precisely, since their statistics are given in figures and not

<sup>7</sup>[https://github.com/ninalu/urlearning-cpp/tree/master/triplet\\_data](https://github.com/ninalu/urlearning-cpp/tree/master/triplet_data)

exactly in tables, though judging from their figures it appears that GRaSP<sub>2</sub> is dominating for adjacency precision and recall, arrowhead recall, and most results for arrowhead precision. Timing results are not reported by Lu et al.; we include these to show that GRaSP<sub>2</sub> returns quickly for all of these examples, where we know (personal communication) that some of the results for Triple A\* take much longer. Adjacencies in these graphs are sampled with uniform probability, “Edge-prob”.

Edge-prob	0.03	0.04	0.05	0.06	0.07	0.08
Precision	0.964	0.976	0.979	0.980	0.982	0.976
Recall	0.985	0.982	0.986	0.986	0.985	0.985
F1	0.974	0.979	0.983	0.983	0.983	0.980

Table 1: GRaSP<sub>2</sub> Adjacency Statistics

Edge-prob	0.03	0.04	0.05	0.06	0.07	0.08
Precision	0.907	0.914	0.933	0.949	0.946	0.945
Recall	0.897	0.916	0.933	0.952	0.952	0.955
F1	0.898	0.913	0.932	0.950	0.948	0.950

Table 2: GRaSP<sub>2</sub> Arrowhead Statistics

Edge-prob	0.03	0.04	0.05	0.06	0.07	0.08
Seconds	0.405	0.755	1.403	2.703	4.795	7.161

Table 3: GRaSP<sub>2</sub> Timing Statistics

## F.2 AIRFOIL EXAMPLE

Figure 4 gives the results of running GRaSP<sub>2</sub>, PC, and fGES on the Airfoil empirical example described in Section 6. GRaSP<sub>2</sub> gets the same uniquely frugal result as SP. To improve readability, we use the names of the variables (instead of numerals) to label the vertices.

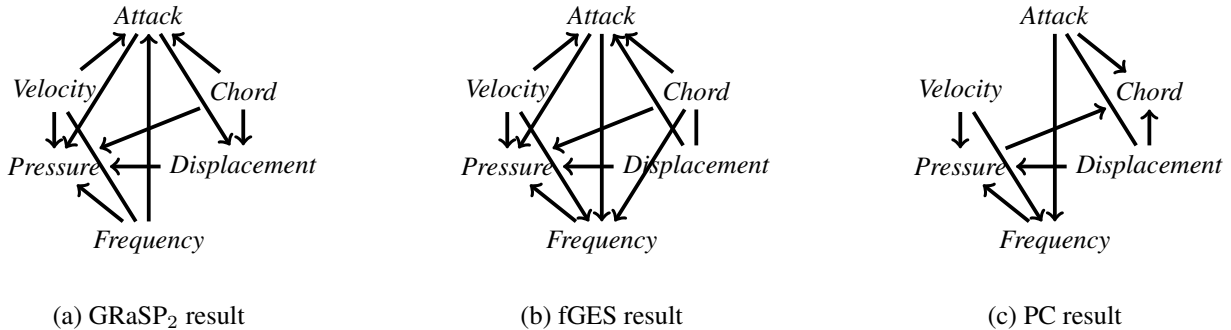


Figure 4: Results of algorithms on NASA airfoil experiment.

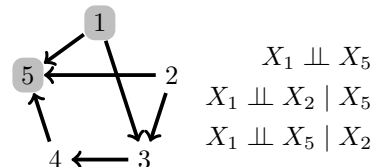
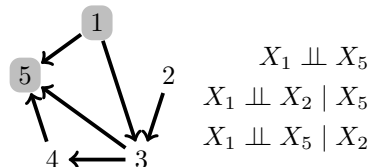
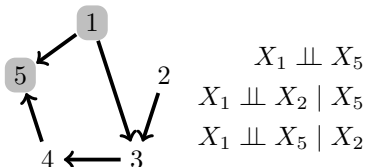
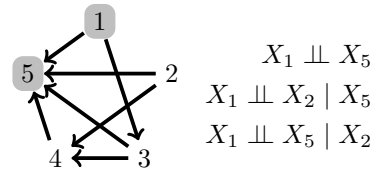
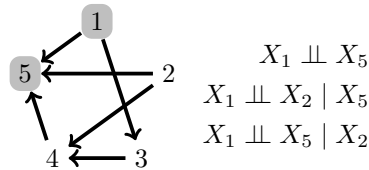
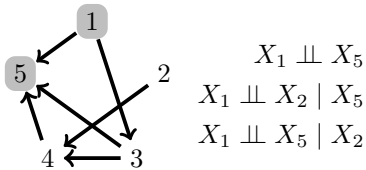
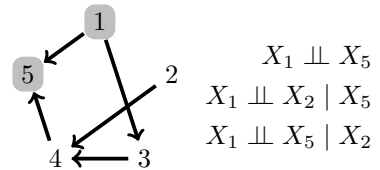
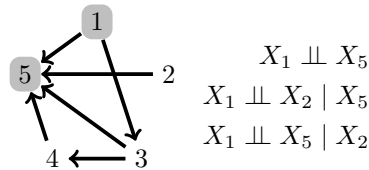
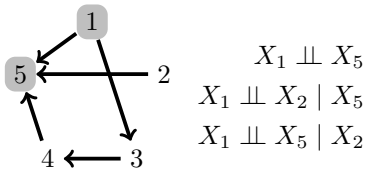
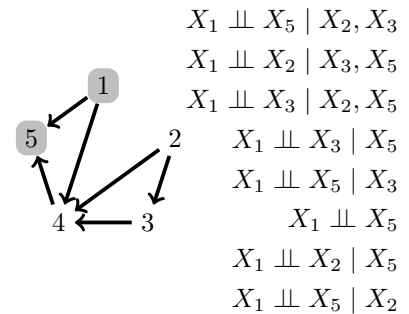
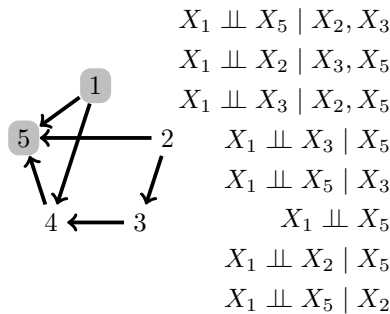
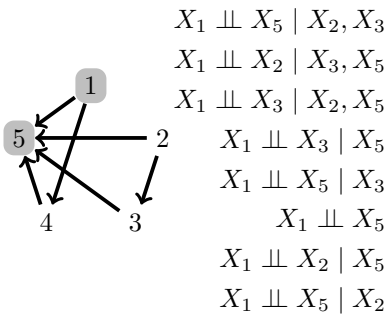
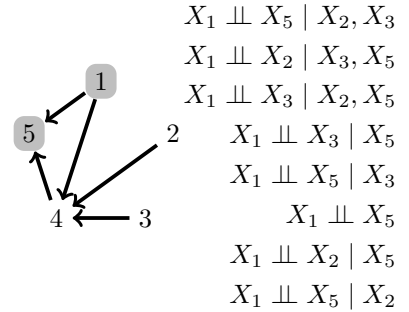
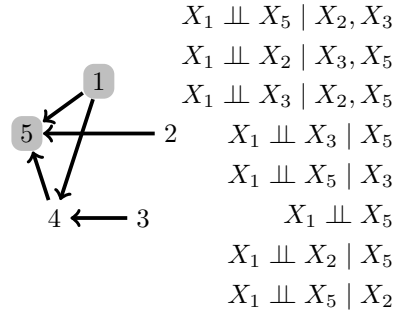
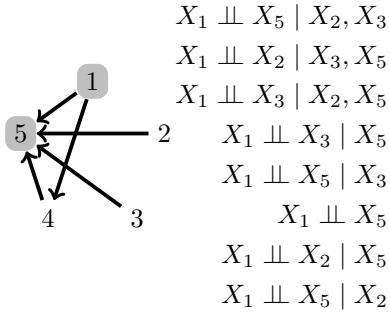
Note that both the GRaSP<sub>2</sub> and FGES results use the linear, Gaussian BIC score with a penalty multiplier of 2. For the GRaSP<sub>2</sub> result in (a), *Attack* is not exogenous, which is counter-intuitive, since it is experimentally controlled. Allowing for latent variables could resolve this issues. However, we leave the development of such an algorithm to future work. On the other hand, the FGES result in (b) is notably not the same as the SP result and so is not frugal. Also, the orientation between *Attack* and *Displacement* is reversed.

The PC result in (c), which uses the zero partial correlation test with a significance level of 0.01, in fact has fewer edges than the frugal result and makes *Chord*, another experimental variable, endogenous. Causally, PC is giving incorrect and incomplete information.

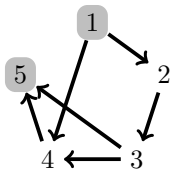


## G UNIT TESTS

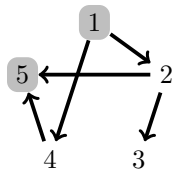
We consider path cancellations in DAGs between pairs of vertices, one of which is exogenous, connected by two or more unique treks. Furthermore, the path cancellations we consider elicit a marginal independence between the two vertices in question. Below, we enumerate all possible path cancellations of this type (up to vertex relabeling). Each graph illustrates a case where an unfaithful marginal independence is elicited between the two gray vertices due to path cancellation. A complete list of all unfaithful CI relations (symmetry assumed) where the independent sets are singletons is also provided for each graph.



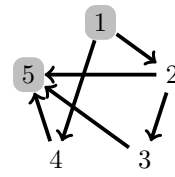




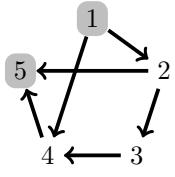
$X_1 \perp\!\!\!\perp X_5$



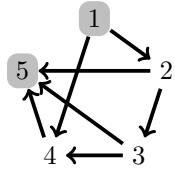
$X_1 \perp\!\!\!\perp X_5$



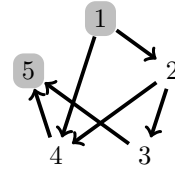
$X_1 \perp\!\!\!\perp X_5$



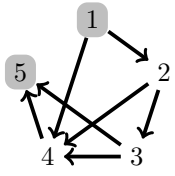
$X_1 \perp\!\!\!\perp X_5$



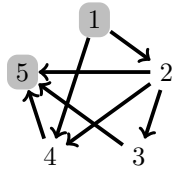
$X_1 \perp\!\!\!\perp X_5$



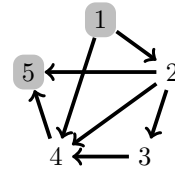
$X_1 \perp\!\!\!\perp X_5$



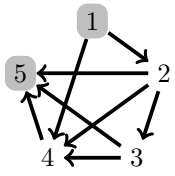
$X_1 \perp\!\!\!\perp X_5$



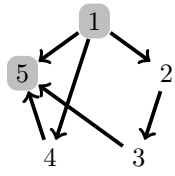
$X_1 \perp\!\!\!\perp X_5$



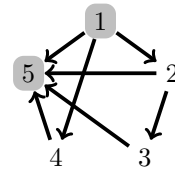
$X_1 \perp\!\!\!\perp X_5$



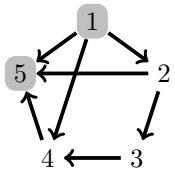
$X_1 \perp\!\!\!\perp X_5$



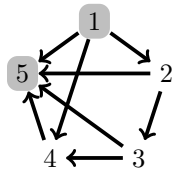
$X_1 \perp\!\!\!\perp X_5$



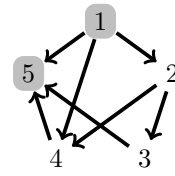
$X_1 \perp\!\!\!\perp X_5$



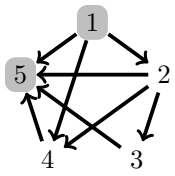
$X_1 \perp\!\!\!\perp X_5$



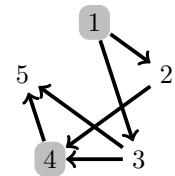
$X_1 \perp\!\!\!\perp X_5$



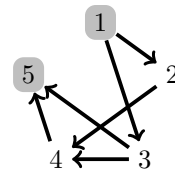
$X_1 \perp\!\!\!\perp X_5$



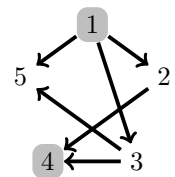
$X_1 \perp\!\!\!\perp X_5$



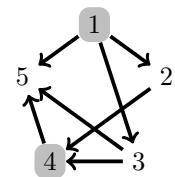
$X_1 \perp\!\!\!\perp X_4$



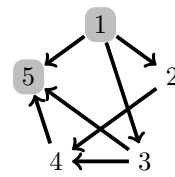
$X_1 \perp\!\!\!\perp X_5$



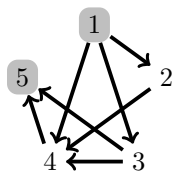
$X_1 \perp\!\!\!\perp X_4$



$X_1 \perp\!\!\!\perp X_4$



$X_1 \perp\!\!\!\perp X_5$



$X_1 \perp\!\!\!\perp X_5$

## References

- David Maxwell Chickering. A transformational characterization of equivalent Bayesian network structures. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, page 87–98, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc.
- David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of machine learning research*, 3: 507–554, 2002.
- Malcolm Forster, Garvesh Raskutti, Reuben Stern, and Naftali Weinberger. The frugal inference of causal relations. *The British Journal for the Philosophy of Science*, 2020.
- Robert E Kass and Paul W Vos. *Geometrical Foundations of Asymptotic Inference*. John Wiley & Sons, Hoboken, NJ, 2011.
- Ni Y Lu, Kun Zhang, and Changhe Yuan. Improving causal discovery by optimal Bayesian network learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 8741–8748, 2021.
- Fatemeh Mohammadi, Caroline Uhler, Charles Wang, and Josephine Yu. Generalized permutohedra from probabilistic graphical models. *SIAM Journal on Discrete Mathematics*, 32:64–93, 2018.
- Judea Pearl. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann Publishers Inc., 1988.
- Joseph Ramsey, Jiji Zhang, and Peter Spirtes. Adjacency-faithfulness and conservative causal inference. In *Proceedings of the conference on Uncertainty in artificial intelligence*, pages 401–408, 2006.
- Garvesh Raskutti and Caroline Uhler. Learning directed acyclic graph models based on sparsest permutations. *Stat*, 7, 2018.
- Liam Solus, Yuhao Wang, and Caroline Uhler. Consistency guarantees for greedy permutation-based causal inference algorithms. *Biometrika*, 108:795–814, 2021.
- Wolfgang Spohn. On the properties of conditional independence. In Paul Humphreys, editor, *Patrick Suppes: Scientific Philosopher: Volume 1. Probability and Probabilistic Causality*, pages 173–196. Springer Netherlands, 1994.
- Milan Studený. *Probabilistic Conditional Independence Structures*. Information Science and Statistics. Springer, 2005.
- Thomas Verma and Judea Pearl. Causal networks: semantics and expressiveness. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, pages 69–78, 1988.
- Jiji Zhang. A comparison of three Occam’s razors for Markovian causal models. *The British journal for the philosophy of science*, 64:423–448, 2013.
- Jiji Zhang and Peter Spirtes. Detection of unfaithfulness and robust causal inference. *Minds and Machines*, 18(2):239–271, 2008.