
Necessary and sufficient graphical conditions for optimal adjustment sets in causal graphical models with hidden variables

Jakob Runge

German Aerospace Center

Institute of Data Science

07745 Jena, Germany

and

Technische Universität Berlin

10623 Berlin, Germany

`jakob.runge@dlr.de`

Abstract

The problem of selecting optimal backdoor adjustment sets to estimate causal effects in graphical models with hidden and conditioned variables is addressed. Previous work has defined optimality as achieving the smallest asymptotic estimation variance and derived an optimal set for the case without hidden variables. For the case with hidden variables there can be settings where no optimal set exists and currently only a sufficient graphical optimality criterion of limited applicability has been derived. In the present work optimality is characterized as maximizing a certain adjustment information which allows to derive a necessary and sufficient graphical criterion for the existence of an optimal adjustment set and a definition and algorithm to construct it. Further, the optimal set is valid if and only if a valid adjustment set exists and has higher (or equal) adjustment information than the Adjust-set proposed in Perković et al. [Journal of Machine Learning Research, 18: 1–62, 2018] for any graph. The results translate to minimal asymptotic estimation variance for a class of estimators whose asymptotic variance follows a certain information-theoretic relation. Numerical experiments indicate that the asymptotic results also hold for relatively small sample sizes and that the optimal adjustment set or minimized variants thereof often yield better variance also beyond that estimator class. Surprisingly, among the randomly created setups more than 90% fulfill the optimality conditions indicating that also in many real-world scenarios graphical optimality may hold.

1 Introduction

A standard problem setting in causal inference is to estimate the causal effect between two variables given a causal graphical model that specifies qualitative causal relations among observed variables [Pearl, 2009], including a possible presence of hidden confounding variables. The graphical model then allows to employ graphical criteria to identify valid adjustment sets, the most well-known being the *backdoor criterion* [Pearl, 1993] and the *generalized adjustment criterion* [Shpitser et al., 2010, Perković et al., 2015, 2018], providing a complete identification of all valid adjustment sets. Estimators of causal effects based on such a valid adjustment set as a covariate are then unbiased, but for different adjustment sets the estimation variance may strongly vary. An *optimal adjustment set* may be characterized as one that has minimal asymptotic estimation variance. In **current work**, following Kuroki and Cai [2004] and Kuroki and Miyakawa [2003], Henckel et al. [2019] (abbreviated

HPM19 in the following) showed that graphical optimality always holds for linear models in the causally sufficient case where all relevant variables are observed. In Witte et al. [2020] an alternative characterization of the optimal adjustment set is discussed and the approach was integrated into the IDA algorithm [Maathuis et al., 2009, 2010] that does not require the causal graph to be known. Rotnitzky and Smucler [2019] extended the results in HPM19 to asymptotically linear non-parametric graphical models. HPM19’s optimal adjustment set holds for the causally sufficient case (no hidden variables) and the authors gave an example with hidden variables where optimality does not hold in general, i.e., the optimal adjustment set depends on the coefficients and noise terms (more generally, the distribution), rather than just the graph. Most recently, Smucler et al. [2021] (SSR20) partially extended these results to the non-parametric hidden variables case together with *dynamic treatment regimes*, i.e., conditional causal effects. SSR20 provide a sufficient criterion for an optimal set to exist and a definition based on a certain undirected graph-construction using a result by van der Zander et al. [2019]. However, their sufficient criterion is very restrictive and a current major open problem is a *necessary* and sufficient condition for an optimal adjustment set to exist in the hidden variable case and a corresponding definition of an optimal set.

My **main theoretical contribution** is a solution to this problem. Optimality for conditional causal effects in the hidden variables case is fully characterized by an information-theoretic approach involving a certain difference of conditional mutual informations among the observed variables termed the adjustment information. Maximizing the adjustment information formalizes the common intuition to choose adjustment sets that maximally constrain the effect variable and minimally constrain the cause variable. This allows to derive a necessary and sufficient graphical criterion for the existence of an optimal adjustment set. The derived optimal adjustment set also has the property of minimum cardinality, i.e., no node can be removed without sacrificing optimality. Further, the optimal set is valid if and only if a valid adjustment set exists and has higher (or equal) adjustment information than the Adjust-set proposed in Perković et al. [2018] for any graph, whether graphical optimality holds or not. The results translate to minimal asymptotic estimation variance for a class of estimators whose asymptotic variance follows a certain information-theoretic relation that, at present, I could only verify theoretically for the linear case. As **practical contributions** the paper provides extensive numerical experiments that corroborate the theoretical results and show that the optimal adjustment set or minimized variants thereof often yield better variance also beyond the theoretically analyzed estimator class. Code is available in the python package <https://github.com/jakobrunge/tigramite>. More detailed preliminaries, proofs, algorithms, and further numerical experiments are given in the Supplementary Material.

1.1 Preliminaries and problem setting

We consider causal effects in causal graphical models over a set of variables \mathbf{V} with a joint distribution $\mathcal{P} = \mathcal{P}(\mathbf{V})$ that is consistent with an acyclic directed mixed graph (ADMG) $\mathcal{G} = (\mathbf{V}, \mathcal{E})$. Two nodes can have possibly more than one edge which can be *directed* (\leftarrow) or *bi-directed* (\leftrightarrow). See Fig. 1A for an example. Kinships are defined as usual: parents $pa(X)$ for “ $\bullet \rightarrow X$ ”, spouses $sp(X)$ for “ $X \leftrightarrow \bullet$ ”, children $ch(X)$ for “ $X \rightarrow \bullet$ ”. These sets all exclude X . Correspondingly descendants $des(X)$ and ancestors $an(X)$ are defined, which, on the other hand, both include X . The mediator nodes on causal paths from X to Y are denoted $\mathbf{M} = \mathbf{M}(X, Y)$ and exclude X and Y . For detailed preliminaries, including the definition of open and blocked paths, see Supplementary Section A. In this work we only consider a univariate intervention variable X and effect variable Y . We simplify set notation and denote unions of variables as $\{W\} \cup \mathbf{M} \cup \mathbf{A} = \mathbf{WMA}$.

A (possibly empty) set of adjustment variables \mathbf{Z} for the total causal effect of X on Y in an ADMG is called *valid* relative to (X, Y) if the interventional distribution for setting $do(X = \mathbf{x})$ [Pearl, 2009] factorizes as $p(Y|do(X = \mathbf{x})) = \int p(Y|\mathbf{x}, \mathbf{z})p(\mathbf{z})d\mathbf{z}$ for non-empty \mathbf{Z} and as $p(Y|do(X = \mathbf{x})) = p(Y|\mathbf{x})$ for empty \mathbf{Z} . Valid adjustment sets, the set of which is here denoted \mathcal{Z} , can be read off from a given causal graph using the generalized adjustment criterion [Perković et al., 2015, 2018] which generalizes Pearl’s back-door criterion [Pearl, 2009]. To this end define

$$\mathbf{forb}(X, Y) = X \cup des(Y\mathbf{M}) \tag{1}$$

(henceforth just denoted as **forb**). A set \mathbf{Z} is valid if both of the following conditions hold: (i) $\mathbf{Z} \cap \mathbf{forb} = \emptyset$, and (ii) all non-causal paths from X to Y are blocked by \mathbf{Z} . An adjustment set is called *minimal* if no strict subset of \mathbf{Z} is still valid. The validity conditions can in principle be manually checked directly from the graph, but, more conveniently, Perković et al. [2018] define an

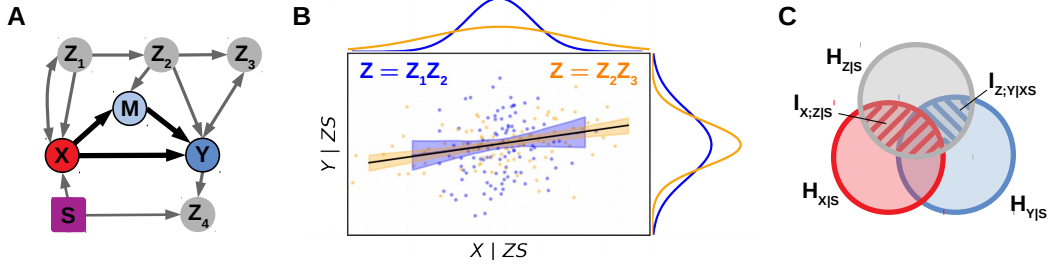


Figure 1: (A) Problem setting of optimal adjustment sets in causal graphs with hidden variables represented through bi-directed edges. The goal is to estimate the total causal effect of X on Y potentially through mediators M , and given conditioned variables S . The task is to select a valid adjustment set Z such that the estimator has minimal asymptotic variance. (B) illustrates two causal effect estimates for a linear Gaussian model consistent with the graph in A, see discussion in text. (C) For a certain class of estimators a minimal asymptotic estimation variance can be translated into an information-theoretical optimization problem, here visualized in a Venn diagram. An optimal adjustment set Z must maximize the adjustment information $I_{Z;Y|XS} - I_{X;Z|S}$ (blue and red hatched, respectively).

adjustment set called ‘Adjust’ that is valid if and only if a valid adjustment set exist. In our setting including conditioning variables S we call this set the *valid ancestors* defined as

$$\mathbf{vancs}(X, Y, S) = \mathbf{an}(XYS) \setminus \mathbf{forb} \quad (2)$$

and refer to this set as \mathbf{vancs} or Adjust-set.

Our quantity of interest is the average total causal effect of an intervention to set X to x vs. x' on the effect variable Y given a set of selected (conditioned) variables $S = s$

$$\Delta_{yxx'|s} = E(Y|do(x), s) - E(Y|do(x'), s). \quad (3)$$

We denote an estimator given a valid adjustment set Z as $\hat{\Delta}_{yxx'|s,z}$. In the linear case $\Delta_{yxx'|s}$ for $x = x' + 1$ corresponds to the regression coefficient $\beta_{YX \cdot ZS}$ in the regression of Y on X , Z , and S . The ordinary least squares (OLS) estimator $\hat{\beta}_{YX \cdot ZS}$ is a consistent estimator of $\beta_{YX \cdot ZS}$.

Figure 1A illustrates the **problem setting**: We are interested in the total causal effect of (here univariate) X on Y (conditioned on S), which is here due to a direct link and an indirect causal path through a mediator M . There are six valid backdoor adjustment sets $Z = \{Z_1, Z_2, Z_1Z_2, Z_2Z_3, Z_1Z_3, Z_1Z_2Z_3\}$. $Z_4 \in \mathbf{forb}$ cannot be included in any set because it is a descendant of YM . Here $\mathbf{vancs} = Z_1Z_2S$. All valid adjustment sets remove the bias due to confounding by their definition. The question is which of these valid adjustment sets is statistically optimal in that it minimizes the asymptotic estimation variance? More formally, the task is, given a graph \mathcal{G} and (X, Y, S) , to chose a valid optimal set $Z_{\text{optimal}} \in \mathcal{Z}$ such that the causal effect estimator’s asymptotic variance $\text{Var}(\hat{\Delta}_{yxx'|s,z}) = E[(\Delta_{yxx'|s} - \hat{\Delta}_{yxx'|s,z})^2]$ is minimal:

$$Z_{\text{optimal}} \in \text{argmin}_{Z \in \mathcal{Z}} \text{Var}(\hat{\Delta}_{yxx'|s,z}). \quad (4)$$

My proposed approach to optimal adjustment sets is based on information theory [Cover and Thomas, 2006]. The main quantity of interest there is the conditional mutual information (CMI) defined as a difference $I_{X;Y|Z} = H_{Y|Z} - H_{Y|ZX}$ of two (conditional) Shannon entropies $H_{Y|X} = -\int_{x,y} p(x,y) \ln p(y|x) dx dy$. Its main properties are non-negativity, $I_{X;Y|Z} = 0$ if and only if $X \perp\!\!\!\perp Y|Z$, and the chain rule $I_{XW;Y|Z} = I_{X;Y|Z} + I_{W;Y|ZX}$. All random variables in a CMI can be multivariate.

Throughout the present paper we will assume the following.

Assumptions 1 (General setting and assumptions). *We assume a causal graphical model over a set of variables \mathbf{V} with a joint distribution $\mathcal{P} = \mathcal{P}(\mathbf{V})$ that is consistent with an ADMG $\mathcal{G} = (\mathbf{V}, \mathcal{E})$. We assume a non-zero causal effect from X on Y , potentially through a set of mediators M , and given selected conditioned variables S , where $S \cap \mathbf{forb} = \emptyset$. We assume that at least one valid adjustment set (given S) exists and, hence, the causal effect is identifiable (except when stated otherwise). Finally, we assume the usual Causal Markov Condition (implicit in semi-Markovian models) and Faithfulness.*

2 Optimal adjustment sets

2.1 Information-theoretic characterization

Figure 1B illustrates two causal effect estimates for a linear Gaussian model consistent with the graph in Fig. 1A. With $\mathbf{Z} = Z_1Z_2$ (blue) the error is much larger than with $\mathbf{O} = Z_2Z_3$ (orange) for two reasons: \mathbf{Z} constrains the residual variance $Var(Y|\mathbf{Z}\mathbf{S})$ of the effect variable Y less than \mathbf{O} and, on the other hand, \mathbf{Z} constrains the residual variance $Var(X|\mathbf{Z}\mathbf{S})$ of the cause variable X more than \mathbf{O} . Smaller estimator variance also holds for \mathbf{O} compared to any other valid set in \mathcal{Z} here.

We information-theoretically formalize the resulting intuition to choose an adjustment set \mathbf{Z} that maximally constrains the effect variable Y and minimally constrains the cause variable X . In terms of CMI and given selected fixed conditions \mathbf{S} the quantity to maximize can be stated as follows.

Definition 1 (Adjustment information). *Consider a causal effect of X on Y for an adjustment set \mathbf{Z} given a condition set \mathbf{S} . The (conditional) adjustment (set) information, abbreviated $J_{\mathbf{Z}}$, is defined as*

$$J_{XY|\mathbf{S};\mathbf{Z}} \equiv I_{\mathbf{Z};Y|X\mathbf{S}} - I_{X;\mathbf{Z}|\mathbf{S}} \quad (5)$$

$$= \underbrace{H_{Y|X\mathbf{S}} - H_{X|\mathbf{S}}}_{\text{not related to } \mathbf{Z}} - \underbrace{(H_{Y|X\mathbf{Z}\mathbf{S}} - H_{X|\mathbf{Z}\mathbf{S}})}_{\text{adjustment entropy}} \quad (6)$$

$J_{\mathbf{Z}}$ is not necessarily positive if the dependence between X and \mathbf{Z} (given \mathbf{S}) is larger than that between \mathbf{Z} and Y given $X\mathbf{S}$. Equation (6) follows from the CMI definition. Fig. 1C illustrates the two CMI in Eq. (5) in a Venn diagram.

Before discussing the range of estimators for which maximizing the adjustment information $J_{\mathbf{Z}}$ leads to a minimal asymptotic estimation variance in Sect. 2.2, we characterize graphical optimality in an information-theoretic framework. Our goal is to provide graphical criteria for optimal adjustment sets, i.e., criteria that depend only on the structure of the graph \mathcal{G} and not on the distribution.

Definition 2 (Information-theoretical graphical optimality). *Given Assumptions 1 we say that (information-theoretical) graphical optimality holds if there is a $\mathbf{Z} \in \mathcal{Z}$ such that either there is no other $\mathbf{Z}' \neq \mathbf{Z} \in \mathcal{Z}$ or for all other $\mathbf{Z}' \neq \mathbf{Z} \in \mathcal{Z}$ and all distributions \mathcal{P} consistent with \mathcal{G} we have $J_{\mathbf{Z}} \geq J_{\mathbf{Z}'}$.*

My main result builds on the following lemma which relates graphical optimality to information-theoretic inequalities in a necessary and sufficient comparison condition for an optimal set to exist.

Lemma 1 (Necessary and sufficient comparison criterion for existence of an optimal set). *Given Assumptions 1, if and only if there is a $\mathbf{Z} \in \mathcal{Z}$ such that either there is no other $\mathbf{Z}' \neq \mathbf{Z} \in \mathcal{Z}$ or for all other $\mathbf{Z}' \neq \mathbf{Z} \in \mathcal{Z}$ and all distributions \mathcal{P} consistent with \mathcal{G} it holds that*

$$\underbrace{I_{\mathbf{Z}\setminus\mathbf{Z}';Y|\mathbf{Z}'X\mathbf{S}}}_{(i)} \geq \underbrace{I_{\mathbf{Z}'\setminus\mathbf{Z};Y|\mathbf{Z}X\mathbf{S}}}_{(iii)} \quad \text{and} \quad \underbrace{I_{X;\mathbf{Z}'\setminus\mathbf{Z}|\mathbf{Z}\mathbf{S}}}_{(ii)} \geq \underbrace{I_{X;\mathbf{Z}\setminus\mathbf{Z}'|\mathbf{Z}'\mathbf{S}}}_{(iv)}, \quad (7)$$

then graphical optimality holds and \mathbf{Z} is optimal implying $J_{\mathbf{Z}} \geq J_{\mathbf{Z}'}$.

In SSR20 and HPM19 the corresponding conditional independence statements to the terms (iii) and (iv) in the inequalities (7) are used as a sufficient pairwise comparison criterion. However, Lemma 1 shows that for graphical optimality it is not necessary that terms (iii) and (iv) vanish, they just need to fulfill the inequalities (7) for a *necessary* and sufficient criterion.

In principle, Lemma 1 can be used to cross-compare all pairs of sets, but firstly, it is difficult to explicitly evaluate (7) for all distributions \mathcal{P} consistent with \mathcal{G} and, secondly, iterating through all valid adjustment sets is computationally prohibitive even for small graph sizes. As an example, consider a confounding path consisting of 5 nodes. Then this path can be blocked by $2^5 - 1$ different subsets. In the main result of this work (Thm. 3) a necessary and sufficient criterion based purely on graphical properties is given.

2.2 Applicable estimator class

The above characterization only relates optimality of adjustment sets to the adjustment information $J_{\mathbf{Z}}$ defined in Eq. (5), but not to any particular estimator. Now the question is for which class of

causal effect estimators $\widehat{\Delta}_{yxx'|s,z}$ the intuition of maximizing the adjustment information $J_{\mathbf{Z}}$ leads to a minimal asymptotic estimation variance. In its most general form this class is characterized as fulfilling

$$\mathbf{Z}_{\text{optimal}} \in \operatorname{argmax}_{\mathbf{Z} \in \mathcal{Z}} J_{\mathbf{Z}} \Leftrightarrow \operatorname{Var}(\widehat{\Delta}_{yxx'|s,z_{\text{optimal}}}) = \min_{\mathbf{Z} \in \mathcal{Z}} \operatorname{Var}(\widehat{\Delta}_{yxx'|s,z}), \quad (8)$$

where we assume that $\widehat{\Delta}_{yxx'|s,z}$ is consistent due to a valid adjustment set and correct functional model specification. One can also further restrict the class to estimators whose (square-root of the) asymptotic variance can be expressed as

$$\sqrt{\operatorname{Var}(\widehat{\Delta}_{yxx'|s,z})} = f(H_{Y|XZS} - H_{X|ZS}), \quad (9)$$

for a real-valued, strictly monotonously increasing function of the adjustment entropy. Minimizing the adjustment entropy is by Eq. (6) equivalent to maximizing the adjustment information. The following assumption and lemma then relates $J_{\mathbf{Z}} \geq J_{\mathbf{Z}'}$ to the corresponding asymptotic variances of a given estimator.

Assumptions 2 (Estimator class assumption). *The model class of the estimator for the causal effect (3) is correctly specified and its asymptotic variance can be expressed as in relation (9).*

Lemma 2 (Asymptotic variance and adjustment information). *Given Assumptions 1 and an estimator fulfilling Assumptions 2, if and only if for two different adjustment sets $\mathbf{Z}, \mathbf{Z}' \in \mathcal{Z}$ we have $J_{\mathbf{Z}} \geq J_{\mathbf{Z}'}$, then the adjustment set \mathbf{Z} has a smaller or equal asymptotic variance compared to \mathbf{Z}' .*

Proof. By Equations (6) and (9) $J_{\mathbf{Z}} \geq J_{\mathbf{Z}'}$ (for fixed X, Y, \mathbf{S}) is directly related to a smaller or equal asymptotic variance for \mathbf{Z} compared to \mathbf{Z}' , and vice versa. \square

The paper’s theoretical results currently hold for estimators fulfilling relation (9), but at least the main result on graphical optimality in Thm. 3 can also be relaxed to estimators fulfilling the less restrictive relation (8). In this work, we leave the question of which general classes of estimators fulfill either relation (8) or the more restricted relation (9) to further research and only show that it holds for the OLS estimator $\widehat{\beta}_{YX.ZS}$ for Gaussian distributions.

For Gaussians the entropies in (9) are given by $H(Y|XZS) = \frac{1}{2} + \frac{1}{2} \ln(2\pi\sigma_{Y|XZS}^2)$ and $H(X|ZS) = \frac{1}{2} + \frac{1}{2} \ln(2\pi\sigma_{X|ZS}^2)$ where $\sigma(\cdot|\cdot)$ denotes the square-root of the conditional variance. Then

$$\sqrt{\operatorname{Var}(\widehat{\Delta}_{yxx'|s,z})} = \frac{1}{\sqrt{n}} e^{H_{Y|XZS} - H_{X|ZS}} = \frac{1}{\sqrt{n}} \frac{\sigma_{Y|XZS}}{\sigma_{X|ZS}}. \quad (10)$$

This relation is also the basis of the results for the causally sufficient case in Henckel et al. [2019] where it is shown that it holds more generally for causal linear models that do not require the noise terms to be Gaussian.

2.3 Definition of O-set

The optimal adjustment set for the causally sufficient case is simply $\mathbf{P} = pa(YM) \setminus \text{forb}$ and was derived in HPM19 and Rotnitzky and Smucler [2019]. In Section B.2 the derivation is discussed from an information-theoretic perspective. In the case with hidden variables we need to account for bidirected edges “ \leftrightarrow ” which considerably complicate the situation. Then the parents of YM are not sufficient to block all non-causal paths. Further, just like conditioning on parents of YM leads to optimality in the sufficient case since parents constrain information in YM , in the hidden variables case also conditioning on spouses of YM constrains information about YM .

Example A. A simple graph (ADMG) to illustrate this is $X \rightarrow Y \leftrightarrow Z_1$ (shown with an additional \mathbf{S} in Fig. 2A below, or Fig. 4 in SSR20). Here $\mathbf{Z}' = \emptyset = \text{vancs}$ is a valid set, but it is not optimal. Consider $\mathbf{O} = Z_1$, then term (iii) = 0 in the inequalities (7) since $\mathbf{Z}' \setminus \mathbf{O} = \emptyset$. Even though not needed to block non-causal paths (there is none), Z_1 still constrains information in Y while being independent of X (hence, term (iv) = 0) which leads to $J_{\mathbf{O}} > J_{\emptyset}$ according to the inequalities (7).

Not only direct spouses can constrain information in Y as Fig. 2B below illustrates. Since for $W \in YM$ the motif “ $W \leftrightarrow \boxed{C_1} \leftarrow * C_2$ ” (“*” denotes either edge mark) is open, it holds that $I(C_1 C_2; W) = I(C_1; Y) + I(C_2; W|C_1) \geq I(C_1; Y)$ and we can even further increase the first term in the adjustment

information by conditioning also on subsequent spouses. This chain of colliders only ends if we reach a tail or there is no further adjacency. However, we have to make sure that conditioning on colliders does not open non-causal paths. This leads to the notion of a *valid collider path* (related to the notion of a *district* in Evans and Richardson [2014]).

Definition 3 (Valid collider paths). *Given a graph \mathcal{G} , a collider path of W for $k \geq 1$ is defined by a sequence of edges $W \leftrightarrow C_1 \leftrightarrow \dots \leftrightarrow C_k$. We denote the set of path nodes (excluding W) along a path indexed by i as π_W^i . Using the set of valid ancestors $\mathbf{vancs} = \mathbf{an}(XY\mathbf{S}) \setminus \mathbf{forb}$ for the causal effect of X on Y given \mathbf{S} we call a collider path node set π_W^i for $W \in Y\mathbf{M}$ valid wrt. to (X, Y, \mathbf{S}) if for each path node $C \in \pi_W^i$ both of the following conditions are fulfilled:*

$$(1) C \notin \mathbf{forb}, \quad \text{and} \quad (2a) C \in \mathbf{vancs} \text{ or } (2b) C \perp\!\!\!\perp X \mid \mathbf{vancs}. \quad (11)$$

Condition (1) is required for any valid adjustment set. If jointly (2a) and (2b) are not fulfilled, i.e. $C \notin \mathbf{vancs}$ and $C \not\perp\!\!\!\perp X \mid \mathbf{vancs}$, then the collider path stops before C . Our candidate optimal adjustment set is now constructed based on the parents of $Y\mathbf{M}$, valid collider path nodes of $Y\mathbf{M}$, and their parents to ‘close’ these collider paths.

Definition 4 (O-set). *Given Assumptions 1 and the definition of valid colliders in Def. 3, define the set $\mathbf{O}(X, Y, \mathbf{S}) = \mathbf{P} \cup \mathbf{C} \cup \mathbf{P}_{\mathbf{C}}$ where*

$$\mathbf{P} = \mathit{pa}(Y\mathbf{M}) \setminus \mathbf{forb}, \quad \mathbf{C} = \uplus_{W \in Y\mathbf{M}} \uplus_i \{ \pi_W^i : \pi_W^i \text{ is valid wrt. to } (X, Y, \mathbf{S}) \}, \quad \mathbf{P}_{\mathbf{C}} = \mathit{pa}(\mathbf{C}).$$

In the following we will abbreviate $\mathbf{O} = \mathbf{O}(X, Y, \mathbf{S})$. Algorithm C.1 states efficient pseudo-code to construct the \mathbf{O} -set and detect whether a valid adjustment set exists. Since none of the conditions of Def. 3 for adding collider nodes depends on previously added nodes, the algorithm is order-independent. The statement occurring in lines 11 and 21 (“No valid adjustment set exists.”) is proven in Thm. 1. If the graph is a DAG, then lines 4-22 can be omitted. The algorithm is of low complexity and the most time-consuming part is checking for a path in line 12, Def. 3(2b) $C \perp\!\!\!\perp X \mid \mathbf{vancs}$, which can be implemented with (bi-directional) breadth-first search as proposed in van der Zander et al. [2019].

Numerical experiments in Section 3 will show that further interesting adjustment sets are the *minimized* \mathbf{O} -set \mathbf{O}_{\min} , where \mathbf{O} is minimized such that no subset can be removed without making \mathbf{O}_{\min} invalid, and the *collider-minimized* \mathbf{O} -set $\mathbf{O}_{\mathbf{C}\min}$ where only $\mathbf{C}\mathbf{P}_{\mathbf{C}} \setminus \mathbf{P} \subseteq \mathbf{O}$ is minimized such that no collider-subset can be removed without making $\mathbf{O}_{\mathbf{C}\min}$ invalid. Both adjustment sets can be constructed with Alg. C.2 similar to the efficient algorithms in van der Zander et al. [2019]. Also the minimized sets are order-independent since the nodes are removed only after the for-loops. Based on the idea in $\mathbf{O}_{\mathbf{C}\min}$, in the numerical experiments we also consider $\mathit{Adjust}_{X\min}$, where only $\mathit{Adjust} \setminus \mathit{pa}(Y\mathbf{M})$ is minimized and $\mathit{pa}(Y\mathbf{M})$ is always included. Finally, we also evaluate Adjust_{\min} where Adjust is fully minimized.

Before discussing the optimality of the \mathbf{O} -set, we need to assure that it is a valid adjustment set. Similar to the proof given in Perković et al. [2018] for the validity of the \mathbf{vancs} -set (for the case without \mathbf{S}), we can state that the \mathbf{O} -set is valid if and only if a valid adjustment set exists.

Theorem 1 (Validity of \mathbf{O} -set). *Given Assumptions 1 but without a priori assuming that a valid adjustment set exists (apart from the requirement $\mathbf{S} \cap \mathbf{forb} = \emptyset$). If and only if a valid backdoor adjustment set exists, then \mathbf{O} is a valid adjustment set.*

2.4 Graphical optimality

We now move to the question of optimality. It is known that there are graphs where no graphical criterion exists to determine optimality. Examples, discussed later, are the graphs in Figs. 2E,F.

Before stating necessary and sufficient conditions for graphical optimality, I mention that next to the \mathbf{O} -set defined above and the Adjust set \mathbf{vancs} [Perković et al., 2018], I am not aware of any other systematically constructed set that will yield a valid adjustment set for the case with hidden variables. van der Zander et al. [2019] provide algorithms to list all valid adjustment sets, but the question is which of these a user should choose. As mentioned above, Lemma 1 can be used to cross-compare all pairs of sets, but this is not really feasible. Hence, for automated causal effect estimation, rather than the question of whether graphical optimality holds, it is crucial to have a set with better properties than other systematically constructable sets. The following theorem states that

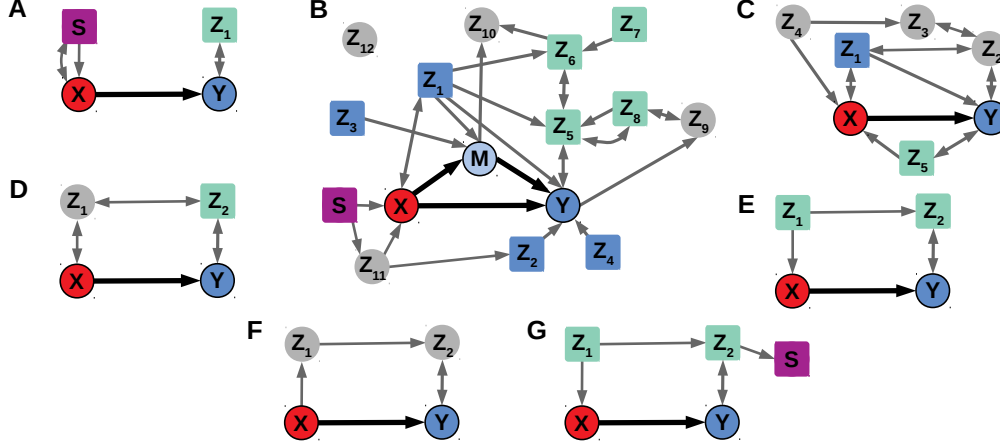


Figure 2: Examples illustrating (optimal) adjustment sets. In all examples the causal effect along causal paths (thick black edges) between X (red circle) and Y (blue circle) potentially through mediators M (light blue circle), and conditioned on some variables S (purple box), is considered. The adjustment set O consists of P (blue boxes) and $CP_C \setminus P$ (green boxes). See main text for details.

the adjustment informations follow $J_O \geq J_{\text{vancs}}$ for any graph (whether graphical optimality holds or not).

Theorem 2 (O-set vs. Adjust-set). *Given Assumptions 1 with O defined in Def. 4 and the Adjust-set defined in Eq. (2), it holds that $J_O \geq J_{\text{vancs}}$ for any graph G . We have $J_O = J_{\text{vancs}}$ only if (1) $O = \text{vancs}$, or (2) $O \subseteq \text{vancs}$ and $X \perp\!\!\!\perp \text{vancs} \setminus O \mid OS$.*

In the following the O -set is illustrated and conditions for graphical optimality are explored. SSR20 provide a sufficient condition for optimality, which states that either all nodes are observed (no bidirected edges exist) or for all observed nodes $V \subset \text{vancs}$. This is a very strict assumption and not fulfilled for any of the examples (except for Example G) discussed in the following.

Example B. Figure 2B depicts a larger example to illustrate the O -set $O = PCP_C$ with $P = Z_1Z_2Z_3Z_4$ (blue boxes) and $CP_C \setminus P = Z_5Z_6Z_7Z_8$ (green boxes). We also have a conditioned variable S . Among P , only Z_1Z_2 are needed to block non-causal paths to X , Z_3Z_4 are only there to constrain information in Y . Here the same holds for the whole set $CP_C \setminus P$ which was constructed from the paths $Y \leftrightarrow Z_5 \leftrightarrow Z_6 \leftarrow Z_7$ and $Y \leftrightarrow Z_5 \leftrightarrow Z_8$ which does not include Z_9 since it is a descendant of YM . Including an independent variable like Z_{12} in O would not decrease the adjustment information J_O , but then O would not be of minimum cardinality anymore (proven in Cor. B.1). Here, again, the condition of SSR20 does not hold (e.g., Z_5 is not an ancestor of XY). O is optimal here which can be seen as follows: For term (iii) in the inequalities (7) to even be non-zero, we would need a valid Z such that $Z \setminus O$ has a path to Y given OSX . But these are all blocked. Note that while Z_{10} or $Z_9 \in Z$ would open a path to Y , both of these are descendants of M or Y and, hence, cannot be in a valid Z . For term (iv) to even be non-zero $O \setminus Z$ would need to have a path to X given ZS . But since any valid Z has to contain Z_1 and Z_2 (or Z_{11}), the only nodes in O with a path to X are parents of YM and paths from these parents to X all need to be blocked for a valid Z . Hence, O is optimal here.

Example C. In Fig. 2C a case is shown where $O = Z_1Z_5$. Z_2 is not part of O because none of the conditions in Def. 3(2) is fulfilled: $Z_2 \notin \text{vancs} = Z_1Z_4Z_5$ and $Z_2 \not\perp\!\!\!\perp X \mid \text{vancs}$. Hence, we call Z_2 an N-node. But Z_2 cannot be part of any valid Z because it has a collider path to X through Z_1 which is always open because it is part of vancs . Hence, term (iii) is always zero. Term (iv) is zero because $O \setminus Z$ is empty for any valid Z here. Here even $J_O > J_Z$ since O is minimal and term (ii) $I_{X;Z \setminus O \mid O} > 0$ for any $Z \neq O$ (generally proven in Corollary B.1).

Example D. The example in Fig. 2D depicts a case with $O = Z_2$ where Z_1 is an N-node. Next to $Z = \emptyset$ another valid set is $Z = Z_1$. Then term (iii) is non-zero and in the same way term (iv) is non-zero. The sufficient pairwise comparison criterion in SSR20 and HPM19 is, hence, not applicable. However, it holds that always term (iii) \leq (i) because the dependence between Z_1 and Y given X is

always smaller than the dependence between Z_2 and Y given X and correspondingly term (iv) \leq (ii). Hence, \mathbf{O} is optimal here. If a link $Y \rightarrow Z_1$ exists, then the only other valid set is $\mathbf{Z} = \emptyset$ and both terms are strictly zero.

Example E. The example in Fig. 2E (Fig. 3 in SSR20 and also discussed in HPM19) is not graphically optimal. Here $\mathbf{O} = Z_1 Z_2$. Other valid adjustment sets are Z_1 or the empty set. From using $Z_1 \perp\!\!\!\perp Y|X$ and $X \perp\!\!\!\perp Z_2|Z_1$ in the inequalities (7) one can derive in information-theoretic terms that both $Z_1 Z_2$ and \emptyset are better than $\mathbf{vancs} = Z_1$, but since $J_{Z_1 Z_2} = J_\emptyset + I_{Z_2;Y|X Z_1} - I_{X;Z_1}$, a superior adjustment set depends on how strong the link $Z_1 \rightarrow X$ vs. $Z_2 \leftrightarrow Y$ is. The graph stays non-optimal also with a link $Z_1 \leftrightarrow Z_2$.

Example F. The example in Fig. 2F is also not graphically optimal. Here $\mathbf{O} = \emptyset$ and Z_2 is an N-node with a non-collider path to X . Other valid adjustment sets are Z_1 and $Z_1 Z_2$. Higher adjustment information here depends on the distribution. Also the same graph with the link $Z_1 \leftrightarrow X$ is non-optimal. If, however, there is another link $Z_1 \rightarrow Y$, then $\mathbf{O} = \emptyset$ is optimal (then Z_1 is a mediator).

Example G. The example in Fig. 2G is only a slight modification of Example E with an added selected condition \mathbf{S} . Then $Z_1, Z_2 \in \mathbf{vancs}$. We still get $\mathbf{O} = Z_1 Z_2$ and this is now optimal since Z_2 is always open and any valid set has to contain Z_1 .

The main result of this work is a set of necessary and sufficient conditions for the existence of graphical optimality and the proof of optimality of the \mathbf{O} -set which is based on the intuition gained in the preceding examples.

Theorem 3 (Necessary and sufficient graphical conditions for optimality and optimality of \mathbf{O} -set). *Given Assumptions 1 and with $\mathbf{O} = \mathbf{PCP}_{\mathbf{C}}$ defined in Def. 4. Denote the set of N-nodes by $\mathbf{N} = sp(\mathbf{YMC}) \setminus (\mathbf{forbOS})$. Finally, given an $N \in \mathbf{N}$ and a collider path $N \leftrightarrow \dots \leftrightarrow C \leftrightarrow \dots \leftrightarrow W$ (including $N \leftrightarrow W$) for $C \in \mathbf{C}$ and $W \in \mathbf{YM}$ (indexed by i) with the collider path nodes denoted by π_i^N (excluding N and W), denote by $\mathbf{O}_{\pi_i^N} = \mathbf{O}(X, Y, \mathbf{S}' = \mathbf{SN}\pi_i^N)$ the \mathbf{O} -set for the causal effect of X on Y given $\mathbf{S}' = \mathbf{S} \cup \{N\} \cup \pi_i^N$. If and only if exactly one valid adjustment set exists, or both of the following conditions are fulfilled, then graphical optimality holds and \mathbf{O} is optimal:*

(I) *For all $N \in \mathbf{N}$ and all its collider paths i to $W \in \mathbf{YM}$ that are inside \mathbf{C} it holds that $\mathbf{O}_{\pi_i^N}$ does not block all non-causal paths from X to Y , i.e., $\mathbf{O}_{\pi_i^N}$ is non-valid,*

and

(II) *for all $E \in \mathbf{O} \setminus \mathbf{P}$ with an open path to X given $\mathbf{SO} \setminus \{E\}$ there is a link $E \leftrightarrow W$ or an extended collider path $E \ast \rightarrow C \leftrightarrow \dots \leftrightarrow W$ inside \mathbf{C} for $W \in \mathbf{YM}$ where all colliders $C \in \mathbf{vancs}$.*

Condition (I) and (II) essentially rule out the two canonical cases in Examples F and E, respectively, on which non-optimality in any graph is based. Applied to the examples, we obtain that in Example A Cond. (I) holds since no N-node exists and Cond. (II) holds since $X \perp\!\!\!\perp Z_1 | S$. In Example B also no N-node exists and Cond. (II) holds as $X \perp\!\!\!\perp E | \mathbf{SO} \setminus \{E\}$ for every $E \in \mathbf{O} \setminus \mathbf{P}$. In example C Z_2 is an N-node, but there is a collider path to X through Z_1 which is in \mathbf{vancs} such that Cond. I is fulfilled. Further, while $X \not\perp\!\!\!\perp Z_5 | \mathbf{SO} \setminus \{Z_5\}$, there is a link $Z_5 \leftrightarrow Y$ such that Cond. II holds. In example D Z_1 is an N-node, but it has a bidirected link with X and Cond. (II) holds since $X \perp\!\!\!\perp Z_2 | \mathbf{SO} \setminus \{Z_2\}$. In Example E optimality does not hold, but Cond. (I) actually holds since there is no N-node. Cond. (II) is not fulfilled for $E = Z_1$, which has a path to X given \mathbf{O} and on the extended collider path $Z_1 \rightarrow Z_2 \leftrightarrow Y$ $Z_2 \notin \mathbf{vancs}$. For $\mathbf{Z}' = \emptyset$ and a distribution \mathcal{P}' where the link $Z_2 \leftrightarrow Y$ almost vanishes we then have $J_{\mathbf{O}} < J_{\mathbf{Z}'}$. Example F has an N-node Z_2 and $\mathbf{O}_{\pi_i^N} = \mathbf{O}(X, Y, \mathbf{S}' = Z_2) = Z_1 Z_2$ is valid implying that Cond. (I) does not hold, while Cond. (II) is actually fulfilled with $\mathbf{O} = \emptyset$. For $\mathbf{Z}' = \mathbf{O}_{\pi_i^N} = Z_1 Z_2$ and a distribution \mathcal{P}' where the link $X \rightarrow Z_1$ almost vanishes we then have $J_{\mathbf{O}} < J_{\mathbf{Z}'}$. Example G is optimal since there are no N-nodes and $Z_2 \in \mathbf{vancs}$.

Similar to SSR20, HPM19, and Witte et al. [2020], I also provide results regarding minimality and minimum cardinality for the hidden variables case in the Supplement.

3 Numerical experiments

We now investigate graphical optimality empirically to answer three questions: Firstly, whether for a linear estimator under Assumptions 2 the asymptotically optimal variance also translates into better

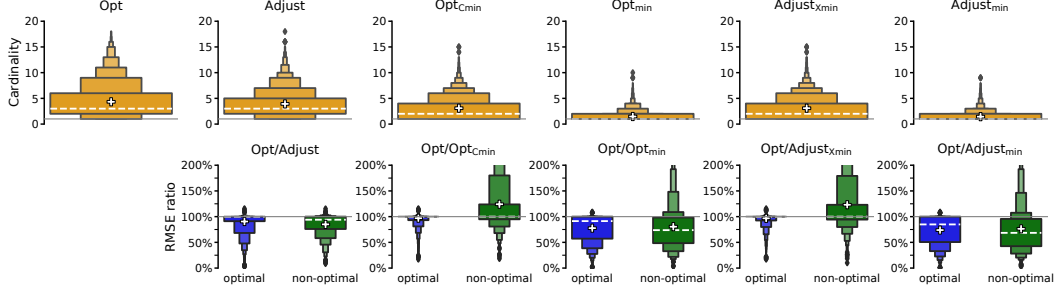


Figure 3: Results of linear experiments with LinReg and sample size $n = 100$. Shown are letter-value plots [Hofmann et al., 2017] of adjustment set cardinalities (top row), as well as RMSE ratios (bottom row) for the \mathbf{O} -set vs. other approaches for optimal configurations (left in blue) and non-optimal configurations (right in green). RMSE was estimated from 100 realizations. The dashed horizontal line denotes the median of the RMSE ratios, and the white ‘plus’ their average. The letter-value plots are interpreted as follows: The widest box shows the 25%–75% range. The next smaller box above (below) shows the 75%–87.5% (12.5%–25%) range and so forth.

finite-sample variance. Secondly, how the \mathbf{O} -set performs in non-optimal settings (according to Thm. 3). Thirdly, how the \mathbf{O} -set and variants thereof perform for estimators not captured by the class for which the theoretical results were derived (Assumptions 2). To this end, we compare the performance of \mathbf{O} , Adjust, \mathbf{O}_{Cmin} , \mathbf{O}_{min} , Adjust $_{Xmin}$, and Adjust $_{min}$ (see definitions in Section 2.3) together with linear least squares estimation (LinReg) on linear models. In the Supplement we also investigate nonlinear models using nearest neighbor regression (kNN), a multilayer perceptron (MLP), random forest regression, and double machine learning for partially linear regression models (DML) [Chernozhukov et al., 2018]. The experiments are based on a generalized additive model and described in detail in Section D. Among these 12,000 randomly created configurations 93% fulfill the optimality conditions in Thm. 3.

The results in Fig. 3 confirm our first hypothesis that for linear experiments with an estimator fulfilling Assumptions 2 and in settings where graphical optimality is fulfilled (Thm. 3) the \mathbf{O} -set either has similar RMSE or significantly outperforms all other tested variants. In particular, \mathbf{O}_{min} and Adjust $_{min}$ are bad choices for this setting. Adjust is intermediate and \mathbf{O}_{Cmin} and Adjust $_{Xmin}$ come closest to \mathbf{O} , but may still yield significantly higher variance.

Secondly, in non-optimal settings (only 7% of configurations) the \mathbf{O} -set still outperforms Adjust (as expected by Thm. 2). Compared to \mathbf{O}_{Cmin} and Adjust $_{Xmin}$ the \mathbf{O} -set leads to worse results for about half of the studied configurations, while \mathbf{O}_{min} and Adjust $_{min}$ are still bad choices. Cardinality is slightly higher for \mathbf{O} compared to all other sets. In Fig. S7 we further differentiate the results by the cardinality of the \mathbf{O} -set and find that for small cardinalities (up to 4) the \mathbf{O} -set has the lowest variance in a majority of cases, but for higher cardinalities either \mathbf{O}_{Cmin} or again \mathbf{O} have the lowest variance (slightly beating Adjust $_{Xmin}$). Hence, either \mathbf{O} or \mathbf{O}_{Cmin} performs best in non-optimal configurations. For very small sample sizes $n = 30$ (see Fig. S2) that become comparable to the adjustment set cardinality, there tends to be a trade-off and smaller cardinality helps. Then \mathbf{O}_{Cmin} tends to be better than \mathbf{O} for high cardinalities also in optimal settings, but here this effect is only present for $n = 30$ and for $n = 50$ already negligible compared to the gain in $J_{\mathbf{O}}$. In Appendix D.2 are RMSE ratios for all combinations of adjustment approaches considered here and it is shown that, in general, results are very similar for other sample sizes.

Thirdly, we investigate non-parametric estimators on linear as well as nonlinear models (implementations described in Section D, results in the figures of Section D.3) The different classes of estimators exhibit quite different behavior. For kNN (Figs. S8,S9) the \mathbf{O} -set has the lowest variance in around 50% of the configurations followed by \mathbf{O}_{Cmin} and \mathbf{O}_{min} . More specifically (Figs. S15,S16), for small \mathbf{O} -set cardinalities up to 2 the \mathbf{O} -set and for higher either \mathbf{O}_{min} or \mathbf{O}_{Cmin} (the latter only in non-optimal configurations) perform best. For nonlinear experiments the results are less clear for \mathbf{O} -set cardinalities greater than 2, but \mathbf{O}_{min} is still a good choice. Regarding RMSE ratios, we see that, for the cases where \mathbf{O} is not the best, the \mathbf{O} -set can have considerably higher variance, while \mathbf{O}_{min} seems to be most robust and may be a better choice if \mathbf{O} is too large. MLP (Figs. S10,S11) behaves much differently. Here in optimal cases neither method outperforms any other for small

O-set cardinalities, but for higher cardinalities (Figs. S15,S16) the **O**-set is best in more than 50% of configurations (slightly less for nonlinear experiments) and the others share the rest (except Adjust_{\min}). For non-optimal cases **O**, $\mathbf{O}_{C_{\min}}$ and $\text{Adjust}_{X_{\min}}$ share the ranks. Regarding RMSE, for linear experiments the **O**-results are almost as optimal as for the LinReg estimator in the optimal setting. However, for non-optimal cases $\mathbf{O}_{C_{\min}}$ can have considerably smaller variance and seems to be a robust option then, similarly to $\text{Adjust}_{X_{\min}}$. Also for nonlinear experiments $\mathbf{O}_{C_{\min}}$ is more robust. The **RF** estimator (Figs. S12,S13) is again different. Here no method clearly is top-ranked, \mathbf{O}_{\min} and Adjust_{\min} are slightly better for linear experiments and **O** for nonlinear experiments. $\mathbf{O}_{C_{\min}}$ and \mathbf{O}_{\min} are more robust regarding RMSE ratios (similar to $\text{Adjust}_{X_{\min}}$). Finally, the **DML** estimator (Fig. S14) was here applied only to linear experiments since its model assumption does not allow for fully nonlinear settings. For optimal settings here **O** is top-ranked in a majority of cases, but closely followed by $\mathbf{O}_{C_{\min}}$ and $\text{Adjust}_{X_{\min}}$. In non-optimal cases for higher **O**-set cardinalities these two seem like a better choice. Quantitatively, $\mathbf{O}_{C_{\min}}$ and $\text{Adjust}_{X_{\min}}$ are the most robust choices.

Overall, the **O**-set and its variants seem to outperform or match the Adjust-variants and whether higher cardinality of the **O**-set reduces performance depends strongly on the estimator and data.

4 Discussion and Conclusions

The proposed adjustment information formalizes the common intuition to choose adjustment sets that maximally constrain the effect variable and minimally constrain the cause variable. The main **theoretical contributions** are a necessary and sufficient graphical criterion for the existence of an optimal adjustment set in the hidden variables case and a definition and algorithm to construct it. To emphasize, graphical optimality implies that the **O**-set is optimal for *any distribution* consistent with the graph. Note that in cases where graphical optimality does not hold, there will still be distributions for which the **O**-set has maximal adjustment information.

Further, the optimal set is valid if and only if a valid adjustment set exists and has smaller (or equal) asymptotic variance compared to the Adjust-set proposed in Perković et al. [2018] for any graph, whether graphical optimality holds or not. This makes the **O**-set a natural choice in automated causal inference analyses. **Practical contributions** comprise Python code to construct adjustment sets and check optimality, as well as extensive numerical experiments that demonstrate that the theoretical results also hold for relatively small sample sizes.

The theoretical **optimality results are limited** to estimators for which the asymptotic variance becomes minimal for adjustment sets with maximal adjustment information (relation (8)). This is fulfilled for least-squares estimators, where even the direct relation (9) holds, but it is unclear whether this also holds for more general classes. The numerical results show that the **O**-set or minimized variants thereof often yield smaller variance also in non-optimal settings and beyond that estimator class. I speculate that further theoretical properties of maximizing adjustment information can be shown because relation (9) for $f(\cdot) = \frac{1}{\sqrt{n}} e^{H_{Y|XZS} - H_{X|ZS}}$ seems related to the lower bound of the estimation variance counterpart to Fano’s inequality (Theorem 8.6.6 in Cover and Thomas [2006]). For estimators sensitive to high-dimensionality one may consider data-driven criteria or penalties to step-wisely minimize the **O**-set. However, estimating, for example, the adjustment information from a potentially small sample size carries considerable errors itself. Another current limitation is that relation (9) only holds for univariate singleton cause variables X . The information-theoretical results, however, also hold for multivariate \mathbf{X} and preliminary results indicate that, while relation (9) does not hold for multivariate \mathbf{X} , the less restrictive relation (8) still seems to hold.

The proposed information-theoretic approach can guide **further research**, for example, to theoretically study relations (8),(9) for other estimators and to address other types of graphs as emerge from the output of causal discovery algorithms and the setting where the graph is unknown [Witte et al., 2020, Maathuis et al., 2009, 2010]. At present, the approach only applies to ADMGs and *Maximal Ancestral Graphs* (MAG) [Richardson and Spirtes, 2002] without selection variables. Last, it remains an open problem to identify optimal adjustment estimands for the hidden variables case based on other criteria such as the front-door formula and Pearl’s general do-calculus [Pearl, 2009].

The results may carry considerable **practical impact** since, surprisingly, among the randomly created configurations more than 90% fulfill the optimality conditions indicating that also in many real-world scenarios graphical optimality may hold. Code is available in the python package <https://github.com/jakobrunge/tigramite>.

Acknowledgments and Disclosure of Funding

I thank Andreas Gerhardus for very helpful comments. This work was funded by the ERC Starting Grant CausalEarth (grant no. 948112).

References

- Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21:C1–C68, 2018.
- Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. John Wiley & Sons, Hoboken, 2006.
- Robin J Evans and Thomas S Richardson. Markovian acyclic directed mixed graphs for discrete data. *The Annals of Statistics*, pages 1452–1482, 2014.
- Leonard Henckel, Emilija Perković, and Marloes H Maathuis. Graphical criteria for efficient total effect estimation via adjustment in causal linear models. *arXiv preprint arXiv:1907.02435*, 2019.
- Heike Hofmann, Hadley Wickham, and Karen Kafadar. Letter-value plots: Boxplots for large data. *Journal of Computational and Graphical Statistics*, 26(3):469–477, 2017.
- Manabu Kuroki and Zhihong Cai. Selection of identifiability criteria for total effects by using path diagrams. In *Proceedings of the 20th conference on Uncertainty in artificial intelligence*, pages 333–340, 2004.
- Manabu Kuroki and Masami Miyakawa. Covariate selection for estimating the causal effect of control plans by using causal diagrams. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 65(1):209–222, 2003.
- Marloes H Maathuis, Markus Kalisch, Peter Bühlmann, et al. Estimating high-dimensional intervention effects from observational data. *The Annals of Statistics*, 37(6A):3133–3164, 2009.
- Marloes H Maathuis, Diego Colombo, Markus Kalisch, and Peter Bühlmann. Predicting causal effects in large-scale systems from observational data. *Nature Methods*, 7:247–8, 2010.
- Judea Pearl. [bayesian analysis in expert systems]: Comment: Graphical models, causality and intervention. *Statistical Science*, 8(3):266–269, 08 1993. doi: 10.1214/ss/1177010894.
- Judea Pearl. *Causality: Models, reasoning, and inference*. Cambridge University Press, 2009.
- Emilija Perković, Johannes Textor, Markus Kalisch, and Marloes H Maathuis. A complete generalized adjustment criterion. In *Uncertainty in Artificial Intelligence-Proceedings of the Thirty-First Conference (2015)*, pages 682–691. AUAI Press, 2015.
- Emilija Perković, Johannes Textor, and Markus Kalisch. Complete graphical characterization and construction of adjustment sets in markov equivalence classes of ancestral graphs. *Journal of Machine Learning Research*, 18:1–62, 2018.
- Thomas Richardson and Peter Spirtes. Ancestral graph markov models. *The Annals of Statistics*, 30(4):962–1030, 08 2002.
- Andrea Rotnitzky and Ezequiel Smucler. Efficient adjustment sets for population average treatment effect estimation in non-parametric causal graphical models. *arXiv preprint arXiv:1912.00306*, 2019.
- Ilya Shpitser, Tyler VanderWeele, and James M Robins. On the validity of covariate adjustment for estimating causal effects. In *Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence*, pages 527–536, 2010.
- E Smucler, F Sapienza, and A Rotnitzky. Efficient adjustment sets in causal graphical models with hidden variables. *Biometrika*, 3 2021. doi: 10.1093/biomet/asab018.

Benito van der Zander, Maciej Liśkiewicz, and Johannes Textor. Separators and adjustment sets in causal graphs: Complete criteria and an algorithmic framework. *Artificial Intelligence*, 270:1–40, 2019.

Janine Witte, Leonard Henckel, Marloes H Maathuis, and Vanessa Didelez. On efficient adjustment in causal graphs. *Journal of Machine Learning Research*, 21(246):1–45, 2020.

Supplementary Material: Necessary and sufficient graphical conditions for optimal adjustment sets in causal graphical models with hidden variables

Jakob Runge

German Aerospace Center
Institute of Data Science
07745 Jena, Germany

and

Technische Universität Berlin
10623 Berlin, Germany
jakob.runge@dlr.de

A Problem setting and preliminaries

A.1 Graph terminology

We consider causal effects in causal graphical models over a set of variables \mathbf{V} with a joint distribution $\mathcal{P} = \mathcal{P}(\mathbf{V})$ that is consistent with an acyclic directed mixed graph (ADMG) $\mathcal{G} = (\mathbf{V}, \mathcal{E})$. Two nodes can have possibly more than one edge which can be *directed* (\leftarrow) or *bi-directed* (\leftrightarrow). We use “*” to denote either edge mark. There can be no loops or directed cycles. See Fig. 1A for an example. The results also hold for *Maximal Ancestral Graphs* (MAG) [Richardson and Spirtes, 2002] without selection variables. A path between two nodes X and Y is a sequence of edges such that every edge occurs only once. A path between X and Y is called *directed or causal* from X to Y if all edges are directed towards Y , else it is called *non-causal*. A node C on a path is called a *collider* if “ $*\rightarrow C \leftarrow *$ ”. Kinships are defined as usual: parents $pa(X, \mathcal{G})$ for “ $\bullet \rightarrow X$ ”, spouses $sp(X, \mathcal{G})$ for “ $X \leftrightarrow \bullet$ ”, children $ch(X, \mathcal{G})$ for “ $X \rightarrow \bullet$ ”, and correspondingly descendants des and ancestors an . We omit the \mathcal{G} in the following since all relations are relative to the graph \mathcal{G} in this paper. Our approach does not involve modified graph constructions as in van der Zander et al. [2019] and other works. A node is an ancestor and descendant of itself, but not a parent/child/spouse of itself. The mediator nodes on causal paths from X to Y are denoted $\mathbf{M} = \mathbf{M}(X, Y)$ and exclude X and Y (different from definitions in other works). For sets of variables the kinship relations correspond to the union of the individual variables. For parent/child/spouse-relationships these exclude the set of variables itself. A path π between X and Y in \mathcal{G} is blocked (or closed) by a node set \mathbf{Z} if (i) π contains a non-collider in \mathbf{Z} or (ii) π contains a collider that is not in $an(\mathbf{Z})$. Otherwise the path π is open (or active/connected) given \mathbf{Z} . Nodes X and Y are said to be m-separated given \mathbf{Z} if every path between them is blocked by \mathbf{Z} , denoted as $X \perp\!\!\!\perp Y | \mathbf{Z}$. In the following we will simplify set notation and denote unions of variables as $\{W\} \cup \mathbf{M} \cup \mathbf{A} = WMA$.

B Further theoretical results and proofs

B.1 Properties of adjustment information

$J_{\mathbf{Z}}$ is not necessarily positive if the dependence between X and \mathbf{Z} (given \mathbf{S}) is larger than that between \mathbf{Z} and Y given $X\mathbf{S}$. By the properties of CMI, it is bounded by

$$-\min(H_{X|\mathbf{S}}, H_{\mathbf{Z}|\mathbf{S}}) \leq J_{XY|\mathbf{S}, \mathbf{Z}} \leq \min(H_{Y|X\mathbf{S}}, H_{\mathbf{Z}|X\mathbf{S}}). \quad (\text{S1})$$

B.2 Causally sufficient case

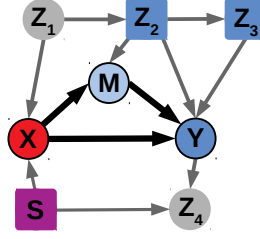


Figure S1: DAG version of graph in Fig. 1A with \mathbf{O} -set shown as blue boxes.

The optimal adjustment set for the causally sufficient case was derived in HPM19 and Rotnitzky and Smucler [2019]. Here the derivation is discussed from an information-theoretic perspective.

Definition B.1 (\mathbf{O} -set in the causally sufficient case). *Given Assumptions 1 restricted to DAGs with no hidden variables, define the set*

$$\mathbf{O} = \mathbf{P} = pa(YM) \setminus \text{forb}.$$

In the causally sufficient case a valid adjustment set always exists and the \mathbf{O} -set is always valid since \mathbf{O} contains no descendants of YM and all non-causal paths from X to Y are blocked since \mathbf{P} blocks all paths from X through parents of YM .

Figure S1 shows an example DAG with a mediator M and conditioned variable S . The \mathbf{O} -set $\mathbf{O} = Z_2Z_3$ is depicted by blue boxes. Compare \mathbf{O} with $\text{vancs} = Z_1Z_2Z_3S$ (Adjust-set in Perković et al. [2018]) in the inequalities (7). Since $Z_1 \perp\!\!\!\perp Y \mid \mathbf{O}XS$, term (iii) is zero and since $\mathbf{O} \setminus \text{vancs} = \emptyset$, also term (iv) is zero. Further, terms (i) and (ii) are both strictly greater than zero (under Faithfulness). Then $J_{\mathbf{O}} > J_{\text{vancs}}$ and under Assumptions 2 by Lemma 2 the \mathbf{O} -set has a smaller asymptotic variance than vancs . Since the parents of YM block all paths from any other valid adjustment sets to Y and because any valid adjustment set \mathbf{Z} has to block paths from X to $pa(YM) \setminus \mathbf{Z}$, $J_{\mathbf{O}} \geq J_{\mathbf{Z}}$ holds in general for any valid set \mathbf{Z} as proven from an information-theoretic perspective in Proposition B.1.

Proposition B.1 (Optimality of \mathbf{O} -set in causally sufficient case). *Given Assumptions 1 restricted to DAGs with no hidden variables and with $\mathbf{O} = \mathbf{P}$ defined in Def. B.1, graphical optimality holds for any graph and \mathbf{O} is optimal.*

Similar to HPM19 and Witte et al. [2020], there also exist results regarding minimality and minimum cardinality which are covered for the hidden variables case in Corollary B.1.

B.3 Hidden variables case

Here we provide some further theoretical results for the general hidden variables case in addition to the lemmas and theorems in the main text.

Corollary B.1 (Minimality and minimum cardinality). *Given Assumptions 1, assume that graphical optimality holds, and, hence, \mathbf{O} is optimal. Further it holds that:*

1. *If \mathbf{O} is not minimal, then $J_{\mathbf{O}} > J_{\mathbf{Z}}$ for all minimal valid $\mathbf{Z} \neq \mathbf{O}$,*
2. *If \mathbf{O} is minimal valid, then \mathbf{O} is the unique set that maximizes the adjustment information $J_{\mathbf{Z}}$ among all minimal valid $\mathbf{Z} \neq \mathbf{O}$,*
3. *\mathbf{O} is of minimum cardinality, that is, there is no subset of \mathbf{O} that is still valid and optimal.*

Another relevant Proposition states that \mathbf{O}_{Cmin} is a subset of vancs , similar to corresponding Lemmas in van der Zander et al. [2019].

Proposition B.2 (Collider-minimized \mathbf{O} -set is a subset of Adjust.). *Given Assumptions 1 with $\mathbf{O} = \text{PCP}_{\mathbf{C}}$ defined in Def. 4 and the \mathbf{O}_{Cmin} -set constructed with Alg. C.2 it holds that $\mathbf{O}_{\text{Cmin}} \subseteq \text{vancs}$.*

B.4 Proof of Lemma 1

Lemma (Necessary and sufficient comparison criterion for existence of an optimal set). Given Assumptions 1, if and only if there is a $\mathbf{Z} \in \mathcal{Z}$ such that either there is no other $\mathbf{Z}' \neq \mathbf{Z} \in \mathcal{Z}$ or for all other $\mathbf{Z}' \neq \mathbf{Z} \in \mathcal{Z}$ and all distributions \mathcal{P} consistent with \mathcal{G} it holds that

$$\begin{aligned} \underbrace{I_{\mathbf{Z} \setminus \mathbf{Z}'; Y | \mathbf{Z}' X S}}_{(i)} &\geq \underbrace{I_{\mathbf{Z}' \setminus \mathbf{Z}; Y | \mathbf{Z} X S}}_{(iii)}, \quad \text{and} \\ \underbrace{I_{X; \mathbf{Z}' \setminus \mathbf{Z} | \mathbf{Z} S}}_{(ii)} &\geq \underbrace{I_{X; \mathbf{Z} \setminus \mathbf{Z}' | \mathbf{Z}' S}}_{(iv)}, \end{aligned} \quad (S2)$$

then graphical optimality holds and \mathbf{Z} is optimal implying $J_{\mathbf{Z}} \geq J_{\mathbf{Z}'}$.

Proof. If there is no other \mathbf{Z}' , the statement trivially holds. Assuming there is another \mathbf{Z}' , we prove the two implications as follows by an information-theoretic decomposition.

Define disjoint (possibly empty) sets $\mathbf{R}, \mathbf{B}, \mathbf{A}$ with $\mathbf{Z} = \mathbf{A}\mathbf{B}$ and $\mathbf{Z}' = \mathbf{B}\mathbf{R}$ with $\mathbf{B} = \mathbf{Z} \cap \mathbf{Z}'$. Note that if both $\mathbf{R} = \emptyset$ and $\mathbf{A} = \emptyset$, then $\mathbf{Z} = \mathbf{Z}'$. Consider two different ways of applying the chain rule of CMI,

$$\begin{aligned} I_{\mathbf{A}\mathbf{B}\mathbf{R}; Y | X S} - I_{X; \mathbf{A}\mathbf{B}\mathbf{R} | S} \\ = I_{\mathbf{A}\mathbf{B}; Y | X S} + I_{\mathbf{R}; Y | \mathbf{A}\mathbf{B} X S} - I_{X; \mathbf{A}\mathbf{B} | S} - I_{X; \mathbf{R} | \mathbf{A}\mathbf{B} S} \end{aligned} \quad (S3)$$

$$= I_{\mathbf{B}\mathbf{R}; Y | X S} + I_{\mathbf{A}; Y | \mathbf{B}\mathbf{R} X S} - I_{X; \mathbf{B}\mathbf{R} | S} - I_{X; \mathbf{A} | \mathbf{B}\mathbf{R} S}, \quad (S4)$$

from which with $J_{\mathbf{Z}} = I_{\mathbf{A}\mathbf{B}; Y | X S} - I_{X; \mathbf{A}\mathbf{B} | S}$ and $J_{\mathbf{Z}'} = I_{\mathbf{B}\mathbf{R}; Y | X S} - I_{X; \mathbf{B}\mathbf{R} | S}$ it follows that

$$\begin{aligned} J_{\mathbf{Z}} - J_{\mathbf{Z}'} \\ + \underbrace{I_{\mathbf{A}; Y | \mathbf{B}\mathbf{R} X S}}_{(i)} + \underbrace{I_{X; \mathbf{R} | \mathbf{A}\mathbf{B} S}}_{(ii)} - \underbrace{I_{\mathbf{R}; Y | \mathbf{A}\mathbf{B} X S}}_{(iii)} - \underbrace{I_{X; \mathbf{A} | \mathbf{B}\mathbf{R} S}}_{(iv)}. \end{aligned} \quad (S5)$$

The inequalities (S2) then read

$$\begin{aligned} \underbrace{I_{\mathbf{A}; Y | \mathbf{B}\mathbf{R} X S}}_{(i)} &\geq \underbrace{I_{\mathbf{R}; Y | \mathbf{A}\mathbf{B} X S}}_{(iii)}, \quad \text{and} \\ \underbrace{I_{X; \mathbf{R} | \mathbf{A}\mathbf{B} S}}_{(ii)} &\geq \underbrace{I_{X; \mathbf{A} | \mathbf{B}\mathbf{R} S}}_{(iv)}. \end{aligned} \quad (S6)$$

“if”: If term (i) is greater or equal to term (iii) and term (ii) greater or equal to term (iv), then trivially $J_{\mathbf{Z}} \geq J_{\mathbf{Z}'}$ for all distributions \mathcal{P} .

“only if”: We prove the contraposition that if for all valid \mathbf{Z} there exists a valid $\mathbf{Z}' \neq \mathbf{Z}$ and a distributions \mathcal{P} consistent with \mathcal{G} such that

$$\underbrace{I_{\mathbf{A}; Y | \mathbf{B}\mathbf{R} X S}}_{(i)} < \underbrace{I_{\mathbf{R}; Y | \mathbf{A}\mathbf{B} X S}}_{(iii)}, \quad \text{or} \quad \underbrace{I_{X; \mathbf{R} | \mathbf{A}\mathbf{B} S}}_{(ii)} < \underbrace{I_{X; \mathbf{A} | \mathbf{B}\mathbf{R} S}}_{(iv)}, \quad (S7)$$

then there always exists a modification \mathcal{P}' of the distribution \mathcal{P} such that $J_{\mathbf{Z}} < J_{\mathbf{Z}'}$. This is because, in both cases, we can always construct a distribution for which terms (ii) and (i), respectively, become arbitrary close to zero. Consider the two cases as follows:

1) there exists a distribution \mathcal{P} with $I_{\mathbf{A}; Y | \mathbf{B}\mathbf{R} X S} < I_{\mathbf{R}; Y | \mathbf{A}\mathbf{B} X S}$: Since CMIs are always non-negative, it holds that $\mathbf{R} \neq \emptyset$ and there must exist at least one open path between \mathbf{R} and Y where every collider is in $\mathbf{A}\mathbf{B} X S$ and no non-collider is in $\mathbf{A}\mathbf{B} X S$. No such open path can pass through X because if X is a non-collider (as for paths continuing on causal paths from X to Y), then the path is blocked, and if X is a collider, then there would be a non-causal path from X to Y given $\mathbf{Z} S$ which would make \mathbf{Z} invalid while \mathbf{Z}' is assumed valid. Correspondingly, no open path from \mathbf{A} (if $\mathbf{A} \neq \emptyset$) to Y given $\mathbf{B}\mathbf{R} X S$, if a path exists at all, can pass through X if \mathbf{Z}' is assumed valid. Now we can construct a distribution \mathcal{P}' with associated structural causal model (SCM) consistent with \mathcal{G} where $I_{\mathbf{A}; Y | \mathbf{B}\mathbf{R} X S} < I_{\mathbf{R}; Y | \mathbf{A}\mathbf{B} X S}$ holds as in \mathcal{P} and still all links “ $U \ast \ast X$ ” for $X \in X$ and $U \notin X M Y$

almost vanish. Consider the three possible links and associated assignment functions in the SCM: (1) “ $X \rightarrow U$ ” with $U := f_U(\dots, X, \dots)$, (2) “ $X \leftarrow U$ ” with $X := f_X(\dots, U, \dots)$, and (3) “ $X \leftrightarrow U$ ” with $X := f_X(\dots, L^U, \dots)$ where L^U denotes one or more latent variables. In each case, to go from \mathcal{P} to \mathcal{P}' , we can modify $f. \rightarrow f'$ where in f' the dependence on the respective argument is replaced by $X \rightarrow cX$, $U \rightarrow cU$, or $L^U \rightarrow cL^U$ for $c \in \mathbb{R}$, and where we consider the limit $c \rightarrow 0$. This modification does not affect $I_{\mathbf{A};Y|\mathbf{BR}XS} < I_{\mathbf{R};Y|\mathbf{AB}XS}$ because the paths contributing to the two CMIs cannot pass through X . On the other hand, then term (ii) $I_{X;\mathbf{R}|\mathbf{ABS}} \rightarrow 0$ because all paths passing through X contain almost zero links and there cannot be a path from \mathbf{R} to X through \mathbf{MY} for a valid \mathbf{Z} . Hence, since in Eq. (S5) term (i) is smaller than term (iii) by assumption, and term (ii) is almost zero, it holds that $J_{\mathbf{Z}} < J_{\mathbf{Z}'}$.

2) there exists a distribution \mathcal{P} with $I_{X;\mathbf{R}|\mathbf{ABS}} < I_{X;\mathbf{A}|\mathbf{BRS}}$: As before, since CMIs are always non-negative, it holds that $\mathbf{A} \neq \emptyset$ and there must exist at least one open path between \mathbf{A} and X where every collider is in \mathbf{BRS} and no non-collider is in \mathbf{BRS} . No such open path can pass through \mathbf{YM} because if any node in \mathbf{YM} is a collider, then the path is blocked, and no path can contain any node in \mathbf{YM} as a non-collider since then either the graph is cyclic or \mathbf{Z}' contains descendants of \mathbf{YM} leading to $\mathbf{Z}' \cap \mathbf{forb} \neq \emptyset$ while \mathbf{Z}' is assumed valid. Correspondingly, no open path from \mathbf{R} (if $\mathbf{R} \neq \emptyset$) to X given \mathbf{ABS} , if a path exists at all, can pass through \mathbf{YM} if \mathbf{Z} is assumed valid. Then, analogous to before, we can construct a \mathcal{P}' with associated SCM consistent with \mathcal{G} where $I_{X;\mathbf{R}|\mathbf{ABS}} < I_{X;\mathbf{A}|\mathbf{BRS}}$ holds and where all links “ $U \ast W$ ” for $W \in \mathbf{YM}$ and $U \notin \mathbf{XMY}$ *almost vanish*. Then term (i) $I_{\mathbf{A};Y|\mathbf{BR}XS} \rightarrow 0$ because all paths contain almost zero links and there cannot be a path from \mathbf{A} to Y where X contains a collider for a valid \mathbf{Z}' since this would constitute a non-causal path. Hence, since in Eq. (S5) term (ii) is smaller than term (iv) by assumption, and term (i) is almost zero, it holds that $J_{\mathbf{Z}} < J_{\mathbf{Z}'}$. \square

B.5 Proof of Proposition B.1

Proposition (Optimality of O-set in causally sufficient case). Given Assumptions 1 restricted to DAGs with no hidden variables and with $\mathbf{O} = \mathbf{P}$ defined in Def. B.1, graphical optimality holds for any graph and \mathbf{O} is optimal.

Proof. The proof is based on Lemma 1 and relation (S5). We will prove that for any DAG \mathcal{G} term (i) \geq (iii) and term (ii) \geq (iv) from which optimality follows by Lemma 1.

We have to show that $I_{\mathbf{A};Y|\mathbf{BR}XS} \geq I_{\mathbf{R};Y|\mathbf{AB}XS}$ and $I_{X;\mathbf{R}|\mathbf{ABS}} \geq I_{X;\mathbf{A}|\mathbf{BRS}}$ where $\mathbf{O} = \mathbf{AB}$ and $\mathbf{Z}' = \mathbf{RB}$ with $\mathbf{B} = \mathbf{O} \cap \mathbf{Z}'$.

Any path from X or $\mathbf{V} \setminus \mathbf{YMOSX}$ to \mathbf{YM} given \mathbf{OS} (denoted by $\boxed{\cdot}$), excluding the causal path from X to Y , features at least one of the following motifs: “ $X, V \ast \ast \boxed{P} \rightarrow W$ ” (excluding “ $X \rightarrow \boxed{P} \rightarrow W$ ”), or “ $V \leftarrow W$ ” where, hence, $V \in \mathbf{forb}$.

Now all paths from a valid adjustment set \mathbf{Z}' with $\mathbf{Z}' \in \mathcal{Z}$ to Y are blocked given \mathbf{OS} : Motif “ $X, V \ast \ast \boxed{P} \rightarrow W$ ” contains a non-collider in \mathbf{OS} and is, hence, blocked. In motif “ $V \leftarrow W$ ” $V \in \mathbf{forb}$. Since $X \notin \mathit{des}(Y)$ (acyclicity) and $\mathbf{Z}' \cap \mathit{des}(Y) = \emptyset$ (validity of \mathbf{Z}'), the paths from \mathbf{Z}' to V either end with a head at V or there must be a collider K that is a descendant of V and hence, $K \in \mathbf{forb}$. Then $K \notin \mathit{an}(\mathbf{OS})$ and $K \notin \mathbf{Z}'$ and the path is therefore blocked. Hence, with $\mathbf{R} \subseteq \mathbf{Z}'$, term (iii) is zero by Markovity.

Term (iv) $I_{X;\mathbf{A}|\mathbf{Z}'S} = 0$ for any valid \mathbf{Z}' because $\mathbf{A} \subseteq \mathit{pa}(\mathbf{YM})$ and then otherwise there would be a non-causal path from X through \mathbf{A} to \mathbf{YM} . \square

B.6 Further Lemmas

Lemma B.1 (Relevant path motifs wrt. the O-set). *Given Assumptions 1 but without a priori assuming that a valid adjustment set exists (apart from the requirement $\mathbf{S} \cap \mathbf{forb} = \emptyset$). With $\mathbf{O} = \mathbf{PCP}_{\mathbf{C}}$ defined in Def. 4 any path from X or $\mathbf{V} \setminus \mathbf{YMOSX}$ to \mathbf{YM} given \mathbf{OS} (denoted by $\boxed{\cdot}$), excluding the causal path from X to Y , features at least one of the following motifs with certain constraints as indicated. We denote $V \in \mathbf{V} \setminus \mathbf{YMOSX}$ and further differentiate nodes in \mathbf{YM} as $W \in \mathbf{YM}$ and in $\mathbf{O} = \mathbf{PCP}_{\mathbf{C}}$ as $C \in \mathbf{C}$ or $P \in \mathbf{P}$ or $P_{\mathbf{C}} \in \mathbf{P}_{\mathbf{C}}$. Last, we denote those collider path nodes not included in the O-set in Alg. C.1 due to not sufficing Def. 3(1) as F with $F \in \mathbf{forb}$ and those not sufficing Def. 3(2a,b) as N with $N \notin \mathbf{forb}$, $N \notin \mathbf{vancs}$, and $N \not\perp X \mid \mathbf{vancs}$:*

- (1a) “ $*-*X \rightarrow \boxed{C} \leftrightarrow$ ”
(1b) “ $*-*X \rightarrow \boxed{P_C} \rightarrow \boxed{C} \leftrightarrow$ ”
(2a) “ $X, V *-* \boxed{P} \rightarrow W$ ” excluding “ $X \rightarrow \boxed{P} \rightarrow W$ ”
(2b) “ $X, V *-* \boxed{P_C} \rightarrow \boxed{C} \leftrightarrow$ ” excluding (1b)
(3a) “ $V \leftarrow W$ ” where, hence, $V \in \mathbf{forb}$
(3b) “ $X, V \leftarrow \boxed{C} \leftrightarrow$ ”
(4a) “ $*-*F \leftrightarrow W$ ” with the constraint $F \notin \mathbf{vancs}$
(4b) “ $*-*F \leftrightarrow \boxed{C} \leftrightarrow$ ” with the constraints $F \notin pa(C)$ and $F \notin \mathbf{vancs}$
(5a) “ $*-*N \leftrightarrow W$ ” with the constraints $N \notin pa(W)$ and $W \notin pa(N)$
(5b) “ $*-*N \leftrightarrow \boxed{C} \leftrightarrow$ ” with the constraint $N \notin pa(C)$

Further it holds that $F, N, X \notin \mathbf{S}$.

Proof. Any path from X or $\mathbf{V} \setminus \mathbf{YMOSX}$ to \mathbf{YM} has to contain a link “ $A*-*B$ ” where $A = X$ or $A \in \mathbf{V} \setminus \mathbf{YMOSX}$ and $B \in \mathbf{YMO}$ where $*-* \in \{\rightarrow, \leftarrow, \leftrightarrow\}$. If we differentiate the left node by X or $V \in \mathbf{V} \setminus \mathbf{YMOSX}$ and the right node by $W \in \mathbf{YM}$ or $C \in \mathbf{C}$ or $P \in \mathbf{P}$ or $P_C \in \mathbf{P}_C$, we can in principle have $2 \cdot 4 \cdot 3 = 24$ link types which are motifs if we consider the adjacent links to A and B . These are listed in the Lemma except for “ $*-*X \rightarrow W$ ” which is part of the causal path from X to Y , “ $X \rightarrow \boxed{P} \rightarrow W$ ” which cannot occur since then $P \in \mathbf{M}$, “ $V \rightarrow W$ ” which cannot occur since \mathbf{P} would contain V or $V \in des(\mathbf{YM})$ leading to a cyclic graph, “ $V \rightarrow C$ ” which cannot occur since \mathbf{P}_C would contain V , and “ $X \leftarrow W$ ” which cannot occur since this implies a cyclic graph.

Regarding the constraints listed in motifs (4a,b) for $F \in \mathbf{forb}$ it holds that $F \notin \mathbf{vancs}$ because $\mathbf{vancs} = an(\mathbf{XYS}) \setminus \mathbf{forb}$ by definition. Further, in (4b) $F \notin pa(C)$ holds because otherwise $C \in \mathbf{forb}$. In motif (5a) $N \notin pa(W)$ holds because $N \notin \mathbf{vancs}$ and $W \notin pa(N)$ holds because $N \notin \mathbf{forb}$. In motif (5b) $N \notin pa(C)$ holds because $C \in \mathbf{vancs}$ contradicts $N \notin \mathbf{vancs}$ and $N \not\perp X \mid \mathbf{vancs}$ with $N \rightarrow C$ contradicts $C \perp X \mid \mathbf{vancs}$. Last, it holds that $F, N, X \notin \mathbf{S}$ because $\mathbf{S} \cap \mathbf{forb} = \emptyset$, $\mathbf{S} \cap X = \emptyset$ by Assumptions 1 and $N \notin \mathbf{vancs}$ while $\mathbf{S} \subseteq \mathbf{vancs}$. \square

Lemma B.2 (Sufficient condition for non-identifiability). *Given Assumptions 1 but without a priori assuming that a valid adjustment set exists (apart from the requirement $\mathbf{S} \cap \mathbf{forb} = \emptyset$). With $\mathbf{O} = \mathbf{PCP}_C$ defined in Def. 4, if on any non-causal path from X to Y given \mathbf{OS} any of the motifs (1a) or (4a) or (4b) for $F = X$ occurs as listed in Lemma B.1, then the causal effect of X on Y (potentially through \mathbf{M}) is not identifiable by backdoor adjustment.*

Proof. If motif (4a) “ $X \leftrightarrow W$ ” for $W \in \mathbf{YM}$ occurs, the case is trivial [Pearl, 2009, Thm. 4.3.1]. In motifs (1a) “ $X \rightarrow \boxed{C} \leftrightarrow$ ” and (4b) “ $X \leftrightarrow \boxed{C} \leftrightarrow$ ” we have that since Def. 3(2b) $C \perp X \mid \mathbf{vancs}$ is not fulfilled, Def. 3(2a) $C \in \mathbf{vancs}$ must be the case. Then every C_k on collider paths to W also fulfills $C_k \in \mathbf{vancs}$ because for all of them $C_k \perp X \mid \mathbf{vancs}$ does not hold since each collider is opened. Hence, there exists a collider path $X * \rightarrow C \leftrightarrow \dots \leftrightarrow W$ where every collider $C \in \mathbf{vancs} = an(\mathbf{XYS}) \setminus \mathbf{forb}$. This path cannot be blocked by any adjustment set (given \mathbf{S}): colliders with $C \in an(\mathbf{S})$ are always open. For colliders with $C \in an(X)$ or $C \in an(Y)$ there is a directed path to X or Y and either this path is open leading to a non-causal path, or an adjustment set contains a non-collider on that directed path which opens the collider C . \square

In Theorem 1 we will prove that the condition in Lemma B.2 is also necessary for non-identifiability by backdoor adjustment. To this end, consider the following Lemmas.

Lemma B.3 (Collider parents fulfill Def. 3). *Given Assumptions 1. With $\mathbf{O} = \mathbf{PCP}_C$ defined in Def. 4, for every $P \in \mathbf{P}_C$ conditions (1), and (2a) or (2b) in Def. 3 hold.*

Proof. Denote a pair $P_C \rightarrow C$ for $C \in \mathbf{C}$ fulfilling conditions (1), and (2a) or (2b) in Def. 3. Firstly, (1) $P_C \notin \mathbf{forb}$ since if $P_C \in des(\mathbf{YM})$ also $C \in des(\mathbf{YM})$ and if $P_C = X$, then by Lemma B.2 no valid adjustment set exists, contrary to Assumptions 1. Secondly, it cannot be that (2a) $P_C \notin \mathbf{vancs}$ and (2b) $P_C \not\perp X \mid \mathbf{vancs}$ because then the path from X to P_C would extend to C and would not be

blocked because $P_C \notin \mathbf{vancs}$. But then also $C \notin \mathbf{vancs}$ and C would not fulfill the conditions in Def. 3. \square

Lemma B.4 (Blockedness of parent-child-motifs). *Given Assumptions 1 with $\mathbf{O} = \mathbf{PCP}_C$ defined in Def. 4. Any path from X or a valid adjustment set \mathbf{Z} with $\mathbf{Z} \in \mathcal{Z}$ to Y containing the motifs (1b), (2a), (2b), (3a), (3b) is blocked given \mathbf{OS} .*

Proof. Motifs (1b), (2a), (2b), and (3b) contain a non-collider in \mathbf{OS} and are, hence, all blocked. In motif (3a) $V \in \mathbf{forb}$. Since $X \notin \mathit{des}(Y)$ (acyclicity) and $\mathbf{Z} \cap \mathit{des}(Y) = \emptyset$ (validity of \mathbf{Z}), the paths from \mathbf{Z} to V either end with a head at V or there must be a collider K that is a descendant of V and hence, $K \in \mathbf{forb}$. Then $K \notin \mathit{an}(\mathbf{OS})$ and $K \notin \mathbf{Z}$ and the path is therefore blocked. \square

Lemma B.5 (Blockedness of F-motifs). *Given Assumptions 1 with $\mathbf{O} = \mathbf{PCP}_C$ defined in Def. 4. Firstly, any path from X to Y containing the motifs (4a) or (4b) for $F \in \mathit{des}(YM)$ is blocked given \mathbf{OS} . Secondly, any path from a valid adjustment set \mathbf{Z} with $\mathbf{Z} \in \mathcal{Z}$ to Y containing the motifs (4a) or (4b) for $F \in \mathit{des}(YM)$ is blocked given $X\mathbf{OS}$.*

Proof. First statement: $F \notin \mathbf{vancs}$ by Lemma B.1 and, hence, in particular $F \notin \mathit{an}(X)$. Then, if a path exists, either the paths from X to F end with a head at F or there must be at least one collider K with $F \in \mathit{an}(K)$ on a path to X . Now $F, K \notin \mathit{an}(\mathbf{OS})$ because $\mathbf{OS} \cap \mathbf{forb} = \emptyset$ and the path is blocked. Secondly, $F \notin \mathit{an}(\mathbf{Z})$ since \mathbf{Z} is valid. Then similarly, if a path exists, either the paths from \mathbf{Z} to F end with a head at F or there must be at least one collider K on a path to \mathbf{Z} with $F \in \mathit{an}(K)$. Now $F, K \notin \mathit{an}(X\mathbf{OS})$ because $\mathbf{OS} \cap \mathbf{forb} = \emptyset$ and $F \notin \mathbf{vancs}$ by Lemma B.1 and the path is blocked. \square

Lemma B.6 (Blockedness of N-motifs). *Given Assumptions 1 with $\mathbf{O} = \mathbf{PCP}_C$ defined in Def. 4. Firstly, any path from X to Y containing the motifs (5a) or (5b) is blocked given \mathbf{OS} . Secondly, any path from a valid adjustment set \mathbf{Z} to Y containing the motifs (5a) or (5b) is blocked given $X\mathbf{OS}$ if \mathbf{Z} does not contain any descendants of N ($\mathbf{Z} \cap \mathit{des}(N) = \emptyset$).*

Proof. First statement: $N \notin \mathbf{vancs}$ by definition of N and, hence, in particular $N \notin \mathit{an}(X)$. Then, if a path exists, either the paths from X to N end with a head at N or there must be at least one collider K with $N \in \mathit{an}(K)$ and $K \notin \mathbf{vancs}$ on a path to X . Now $N, K \notin \mathit{an}(\mathbf{OS})$ can be seen by considering the different parts of \mathbf{O} : $N, K \notin \mathit{an}(\mathbf{PS})$ since $N, K \notin \mathbf{vancs}$ and $N, K \notin \mathit{an}(C)$ for $C \in \mathbf{vancs} \cap \mathbf{CP}_C$. Finally, $N, K \notin \mathit{an}(C)$ for $C \in \mathbf{CP}_C$ with $C \perp\!\!\!\perp X \mid \mathbf{vancs}$ because $N, K \not\perp\!\!\!\perp X \mid \mathbf{vancs}$. Hence, the path is blocked. Second statement: If \mathbf{Z} does not contain any descendants of N , then $N \notin \mathit{an}(\mathbf{Z})$. Then any path from a \mathbf{Z} is blocked by the same reasoning as in the first part with the addition that $N \notin \mathit{an}(X)$ and hence the motif is blocked given $X\mathbf{OS}$. \square

The following Lemma is not needed in this paper, but may be of interest for further research.

Lemma B.7 (Existence of X-N-path). *Given Assumptions 1 with $\mathbf{O} = \mathbf{PCP}_C$ defined in Def. 4. There must be at least one path from X to N (defined in the motifs (5a) or (5b)) that ends with a head at N and where every collider is in \mathbf{vancs} and every non-collider is not in \mathbf{vancs} .*

Proof. By definition of the N-node, $N \not\perp\!\!\!\perp X \mid \mathbf{vancs}$. Now all paths that end with a tail at N are blocked given \mathbf{vancs} because $N \notin \mathit{an}(X)$ and the first collider K coming from N must be blocked because $K \notin \mathbf{vancs}$. Hence, there must be an open path that ends with a head at N and where every collider is in \mathbf{vancs} and every non-collider is not in \mathbf{vancs} as stated. \square

B.7 Proof of Theorem 1

Theorem (Validity of O-set). *Given Assumptions 1 but *without* a priori assuming that a valid adjustment set exists (apart from the requirement $\mathbf{S} \cap \mathbf{forb} = \emptyset$). If and only if a valid backdoor adjustment set exists, then \mathbf{O} is a valid adjustment set.*

Proof. "if": Given that a valid backdoor adjustment set exists, we need to prove that (i) $\mathbf{O} \cap \mathbf{forb} = \emptyset$ with $\mathbf{forb} = X \cup \mathit{des}(YM)$ and (ii) all non-causal paths from X to Y are blocked by \mathbf{O} (given \mathbf{S}). (i) is true by the construction of \mathbf{O} in Def. 4 and Alg. C.1 where nodes $\in \mathit{des}(YM)$ are not added and nodes that are X indicate non-identifiability (see Lemma B.2). By Lemma B.3 also $\mathbf{P}_C \cap \mathit{des}(YM) = \emptyset$ and $X \notin \mathbf{P}_C$ because otherwise no valid adjustment set exists by Lemma B.2.

Lemma B.1 lists all possible motifs on non-causal paths. By Lemma B.2 the occurrence of the motifs (1a) or (4a) or (4b) for $F = X$ renders the effect non-identifiable, contrary to the assumption. Hence

only the remaining motifs can occur. By Lemma B.4 the motifs (1b), (2a), (2b), (3a), (3b) are blocked given \mathbf{OS} . By Lemma B.5 (part one) the motifs (4a,b) for $F \in \text{des}(YM)$ are blocked given \mathbf{OS} . By Lemma B.6 (part one) motifs (5a) and (5b) are blocked given \mathbf{OS} .

“only if” is trivially true since \mathbf{O} is then assumed valid. \square

B.8 Proof of Theorem 2

Theorem (O-set vs Adjust-set). Given Assumptions 1 with \mathbf{O} defined in Def. 4 and the Adjust-set \mathbf{vancs} defined in Eq. (2), it holds that $J_{\mathbf{O}} \geq J_{\mathbf{vancs}}$ for any graph \mathcal{G} . We have $J_{\mathbf{O}} = J_{\mathbf{vancs}}$ only if $\mathbf{O} = \mathbf{vancs}$ or $\mathbf{O} \subseteq \mathbf{vancs}$ and $X \perp\!\!\!\perp \mathbf{vancs} \setminus \mathbf{O} \mid \mathbf{OS}$.

Proof. We directly use the decomposition in Eq. (S5) with $\mathbf{Z} = \mathbf{O} = \mathbf{AB}$ and $\mathbf{Z}' = \mathbf{vancs} = \mathbf{BR}$ with $\mathbf{vancs} = \text{an}(XYS) \setminus \text{forb}$ and the definitions of $\mathbf{R}, \mathbf{B}, \mathbf{A}$ as in Eq. (S5). For term (iii), $I_{\mathbf{R};Y|\mathbf{OXS}}$, to be non-zero, there must be an active path from $\mathbf{R} \subseteq \mathbf{vancs}$ to Y given $X\mathbf{OS}$. By Lemma B.1, Lemma B.4, Lemma B.5 (second part), and Lemma B.6 (second part), the only possibly open motifs on paths from \mathbf{R} to Y given \mathbf{OXS} are “ $\leftarrow N \leftrightarrow W$ ” or “ $\leftarrow N \leftrightarrow \boxed{C} \leftrightarrow$ ” where $\mathbf{R} \cap \text{des}(N) \neq \emptyset$. But since $\mathbf{R} \subseteq \mathbf{vancs}$ and $N \notin \mathbf{vancs}$, \mathbf{R} cannot contain descendants of N . Hence, term (iii) is zero. For term (iv), $I_{X;\mathbf{A}|\mathbf{BRS}} = I_{X;\mathbf{A}|\mathbf{vancs}}$, note that $\mathbf{A} = \mathbf{O} \setminus \mathbf{vancs}$ and, hence, for all $A \in \mathbf{A}$ it holds that $A \perp\!\!\!\perp X \mid \mathbf{vancs}$ since all $A \in \mathbf{A}$ then fulfill Def. 3(2b) (for $A \in \mathbf{P}_{\mathbf{C}}$ see Lemma B.3). Hence, $I_{X;\mathbf{A}|\mathbf{vancs}} = 0$ by Markovity. This proves that $J_{\mathbf{O}} \geq J_{\mathbf{vancs}}$.

We are now left with terms (i) and (ii) in Eq. (S5). By construction of the collider path nodes, $\mathbf{A} \subseteq \mathbf{CP}_{\mathbf{C}}$ is connected to Y (potentially through \mathbf{M}) conditional on $\mathbf{vancs}X$ since \mathbf{vancs} contains all remaining collider nodes in \mathbf{C} . Then by Faithfulness term (i) $I_{\mathbf{A};Y|\mathbf{BRXS}} = I_{\mathbf{A};Y|\mathbf{vancs}X}$ can only be zero if $\mathbf{A} = \emptyset$. Then $\mathbf{O} \subseteq \mathbf{vancs}$. Term (ii), $I_{X;\mathbf{R}|\mathbf{OS}} = 0$ if $\mathbf{R} = \mathbf{vancs} \setminus \mathbf{O} = \emptyset$ or $X \perp\!\!\!\perp \mathbf{vancs} \setminus \mathbf{O} \mid \mathbf{OS}$ together with Faithfulness. \square

B.9 Proof of Proposition B.2

Proposition (Collider-minimized O-set is a subset of Adjust.). Given Assumptions 1 with $\mathbf{O} = \mathbf{PCP}_{\mathbf{C}}$ defined in Def. 4 and the $\mathbf{O}_{\mathbf{Cmin}}$ -set constructed with Alg. C.2 it holds that $\mathbf{O}_{\mathbf{Cmin}} \subseteq \mathbf{vancs}$.

Proof. Define $\mathbf{C}_{\min} = \mathbf{O}_{\mathbf{Cmin}} \setminus \mathbf{P}$. We need to show that $C \in \mathbf{C}_{\min} \Rightarrow C \in \mathbf{vancs}$ for all $C \in \mathbf{O} \setminus \mathbf{P}$. Assume $C \notin \mathbf{vancs}$. Since then all $C \in \mathbf{O} \setminus \mathbf{P}$ fulfill Def. 3(2b) (for $C \in \mathbf{P}_{\mathbf{C}}$ see Lemma B.3), it holds that $C \perp\!\!\!\perp X \mid \mathbf{vancs}$ implying that no link $X \ast \ast C$ exists. If a path exists at all, either (i) there must be at least one collider K with $C \in \text{an}(K)$ and $K \notin \mathbf{vancs}$ on a path to X or (ii) $C \in \text{des}(X)$. We now show that for case (i) C has no open path to X given $\mathbf{SO} \setminus \{C\}$. $K \notin \text{an}(\mathbf{OS})$ can be seen by considering the different parts of \mathbf{OS} : $K \notin \text{an}(\mathbf{PS})$ since $K \notin \mathbf{vancs}$ and $\text{an}(\mathbf{PS}) \subseteq \mathbf{vancs}$. Further, $K \notin \text{an}(\mathbf{vancs} \cap \mathbf{C})$. Finally, $K \notin \text{an}(\mathbf{CP}_{\mathbf{C}} \setminus \mathbf{vancs})$ since $C' \in \mathbf{CP}_{\mathbf{C}} \setminus \mathbf{vancs}$ fulfill (by Def. 3(2b)) $C' \perp\!\!\!\perp X \mid \mathbf{vancs}$ and $K \not\perp\!\!\!\perp X \mid \mathbf{vancs}$. Hence, $X \perp\!\!\!\perp C \mid \mathbf{SO} \setminus \{C\}$ implying that C would be removed in the first loop of Alg. C.2 and $C \notin \mathbf{C}_{\min}$, contrary to assumption.

In case (ii) the directed path from X to C for $C \in \mathbf{C} \setminus \mathbf{P}_{\mathbf{C}}$ is blocked because $\mathbf{P}_{\mathbf{C}} \subseteq \mathbf{O}$ contains all parents of C and $X \notin \mathbf{P}_{\mathbf{C}}$ since we assume identifiability. This implies that C would be removed in the first loop of Alg. C.2 and $C \notin \mathbf{C}_{\min}$, contrary to assumption. Finally, if there exists a directed path from X to $C = P_C \in \mathbf{P}_{\mathbf{C}} \setminus \mathbf{C}$ for $P_C \notin \mathbf{vancs}$ we know that all children $C \in \text{ch}(P_C) \cap \mathbf{CP}$ were removed in the first loop of Alg. C.2. Denote the remaining nodes after the first loop of Alg. C.2 by $\mathbf{O}'_{\mathbf{Cmin}}$. $P_C \notin \mathbf{vancs}$ has no directed path to Y and is separated from Y given $\mathbf{SO}'_{\mathbf{Cmin}}$ because the motif $P_C \rightarrow C \leftrightarrow$ is blocked since $C \notin \text{an}(\mathbf{O}'_{\mathbf{Cmin}})$. This implies that P_C would be removed in the second loop of Alg. C.2 and $P_C \notin \mathbf{C}_{\min}$, contrary to assumption. \square

B.10 Proof of Theorem 3

Theorem (Necessary and sufficient graphical conditions for optimality and optimality of O-set). Given Assumptions 1 and with $\mathbf{O} = \mathbf{PCP}_{\mathbf{C}}$ defined in Def. 4. Denote the set of N-nodes by $\mathbf{N} = \text{sp}(Y\mathbf{MC}) \setminus (\text{forbOS})$. Finally, given an $N \in \mathbf{N}$ and a collider path $N \leftrightarrow \dots \leftrightarrow C \leftrightarrow \dots \leftrightarrow W$ (including $N \leftrightarrow W$) for $C \in \mathbf{C}$ and $W \in Y\mathbf{M}$ (indexed by i) with the collider path nodes denoted by π_i^N (excluding N and W), denote by $\mathbf{O}_{\pi_i^N} = \mathbf{O}(X, Y, \mathbf{S}' = \mathbf{SN}\pi_i^N)$ the O-set for the causal effect of X on Y given $\mathbf{S}' = \mathbf{S} \cup \{N\} \cup \pi_i^N$.

If and only if exactly one valid adjustment set exists, or both of the following conditions are fulfilled, then graphical optimality holds and \mathbf{O} is optimal:

(I) For all $N \in \mathbf{N}$ and all its collider paths i to $W \in \mathbf{YM}$ that are inside \mathbf{C} it holds that $\mathbf{O}_{\pi_i^N}$ does not block all non-causal paths from X to Y , i.e., $\mathbf{O}_{\pi_i^N}$ is non-valid,

and

(II) for all $E \in \mathbf{O} \setminus \mathbf{P}$ with an open path to X given $\mathbf{SO} \setminus \{E\}$ there is a link $E \leftrightarrow W$ or an extended collider path $E \ast \rightarrow C \leftrightarrow \dots \leftrightarrow W$ inside \mathbf{C} for $W \in \mathbf{YM}$ where all colliders $C \in \mathbf{vancs}$.

Proof. If exactly one valid adjustment set exists, then optimality holds by Def. 2 and then this set is \mathbf{O} because \mathbf{O} is always valid if a valid set exists (Lemma 1).

The proof is based on Lemma 1 and relation (S5). We will first prove the “if”-statement by showing that Cond. (I) leads to term (i) \geq (iii) and Cond. (II) leads to term (ii) \geq (iv) from which optimality follows by Lemma 1. Then we prove the “only if”-statement by showing that if either of the two conditions is not fulfilled, then (i) $<$ (iii) or (ii) $<$ (iv) for some distribution \mathcal{P} consistent with \mathcal{G} and graphical optimality does not hold.

“if”: We have to show that if both conditions hold, then $I_{\mathbf{A};Y|\mathbf{B}\mathbf{R}\mathbf{X}\mathbf{S}} \geq I_{\mathbf{R};Y|\mathbf{A}\mathbf{B}\mathbf{X}\mathbf{S}}$ and $I_{X;\mathbf{R}|\mathbf{A}\mathbf{B}\mathbf{S}} \geq I_{X;\mathbf{A}|\mathbf{B}\mathbf{R}\mathbf{S}}$ where $\mathbf{O} = \mathbf{A}\mathbf{B}$ and $\mathbf{Z}' = \mathbf{R}\mathbf{B}$ with $\mathbf{B} = \mathbf{O} \cap \mathbf{Z}'$. Further, we use $\mathbf{A}_{\mathbf{P}} = \mathbf{A} \cap \mathbf{P}$ and $\mathbf{A}_{\mathbf{C}} = (\mathbf{A} \cap \mathbf{C}\mathbf{P}_{\mathbf{C}}) \setminus \mathbf{A}_{\mathbf{P}}$ where $\mathbf{A} = \mathbf{A}_{\mathbf{P}} \cup \mathbf{A}_{\mathbf{C}}$.

Condition (I) directly leads to $I_{\mathbf{A};Y|\mathbf{B}\mathbf{R}\mathbf{X}\mathbf{S}} \geq I_{\mathbf{R};Y|\mathbf{A}\mathbf{B}\mathbf{X}\mathbf{S}}$ as follows.

We subdivide condition (I) into two cases where the former implies the latter: (I.1) There are no N-nodes, i.e., $\mathbf{N} = \emptyset$, or (I.2) for all $N \in \mathbf{N}$ and all its collider paths i it holds that $\mathbf{O}_{\pi_i^N}$ does not block all non-causal paths from X to Y .

If condition (I.1) holds, then there are no N-nodes. If there are no N-motifs on any path from \mathbf{R} to Y , then by Lemma B.1, Lemma B.4, and Lemma B.5 (second part) all paths given $X\mathbf{O}\mathbf{S}$ are blocked and term (iii) is zero by Markovity.

If condition (I.2) holds, then there are N-nodes. By Lemma B.6 (second part) the only possibly open motifs on paths from \mathbf{R} to Y given $\mathbf{O}\mathbf{X}\mathbf{S}$ are “ $\leftarrow N \leftrightarrow W$ ” or “ $\leftarrow N \leftrightarrow \boxed{C} \leftrightarrow$ ” where $\mathbf{R} \cap \mathit{des}(N) \neq \emptyset$. Term (iii), $I_{\mathbf{R};Y|\mathbf{B}\mathbf{X}\mathbf{S}\mathbf{A}} = I_{\mathbf{R};Y|\mathbf{X}\mathbf{S}\mathbf{O}}$, is then always non-zero since, by definition of the N-nodes, there exists at least one collider path $N \leftrightarrow \dots \leftrightarrow \boxed{C} \leftrightarrow \dots \leftrightarrow W$ (including $N \leftrightarrow W$) for $C \in \mathbf{C}$ and $W \in \mathbf{YM}$. To see under which conditions still term (i) \geq (iii) consider two ways of decomposing the following CMI:

$$\begin{aligned} I_{\mathbf{A}\mathbf{R};Y|\mathbf{B}\mathbf{X}\mathbf{S}} &= \underbrace{I_{\mathbf{A};Y|\mathbf{B}\mathbf{X}\mathbf{S}}}_{\text{term (i')}} + \underbrace{I_{\mathbf{R};Y|\mathbf{B}\mathbf{X}\mathbf{S}\mathbf{A}}}_{\text{term (iii)}} \\ &= \underbrace{I_{\mathbf{R};Y|\mathbf{B}\mathbf{X}\mathbf{S}}}_{\text{term (iii')}} + \underbrace{I_{\mathbf{A};Y|\mathbf{B}\mathbf{X}\mathbf{S}\mathbf{R}}}_{\text{term (i)}}. \end{aligned} \quad (\text{S8})$$

From this decomposition we see that term (i) \geq (iii) if and only if term (i') \geq (iii'). Paths from \mathbf{R} to Y via X given $\mathbf{S}\mathbf{X}\mathbf{Z}' \setminus \mathbf{R} = \mathbf{B}\mathbf{S}\mathbf{X}$ are blocked because if X is a collider, then there would be a non-causal path rendering \mathbf{Z}' invalid. Therefore, for term (iii') to be non-zero $\mathbf{Z}'\mathbf{S}$ must contain at least descendants of an N-node N and all its collider path nodes towards W , denoted π_i^N , for at least one path i . Then $\mathbf{R} \cap \mathit{des}(N) \neq \emptyset$ and $\pi_i^N \subseteq \mathbf{B}\mathbf{S}$ such that there exists an open path “ $N \leftrightarrow \boxed{C} \leftrightarrow \dots \leftrightarrow \boxed{C} \leftrightarrow W$ ” (or $N \leftrightarrow W$).

Condition (I.2) now guarantees that for all $N \in \mathbf{N}$ and all collider paths indexed by i the O-set $\mathbf{O}_{\pi_i^N}$, which includes $N\pi_i^N$ as a subset, does *not* block all non-causal paths. By Theorem 1, if $\mathbf{O}_{\pi_i^N}$ is not valid, then no valid adjustment set \mathbf{Z}' containing $N\pi_i^N$ as a subset exists. And this in turn implies that no valid set with $\mathbf{R} \cap \mathit{des}(N) \neq \emptyset$ exists. To show this, assume the contraposition: If there was such a valid set \mathbf{Z}' with $\mathbf{R} \cap \mathit{des}(N) \neq \emptyset$ and $\pi_i^N \subseteq \mathbf{Z}'$, then it would open the collider motif $\ast \rightarrow N \leftrightarrow$ since \mathbf{R} contains descendants of N and lead to an open path “ $N \leftrightarrow \boxed{C} \leftrightarrow \dots \leftrightarrow \boxed{C} \leftrightarrow W$ ” (or $N \leftrightarrow W$). If \mathbf{Z}' is still valid, it must block all paths from X that end with an arrowhead at N . But then also $\mathbf{Z}' \cup \{N\}$ is valid. Note that since $N \notin \mathbf{forb}$, $\mathbf{S} \cap \mathbf{forb} = \emptyset$, and $\pi_i^N \cap \mathbf{forb} = \emptyset$ since

$\pi_i^N \subseteq \mathbf{C}$, the validity of $\mathbf{O}_{\pi_i^N}$ depends only on its ability to block non-causal paths. Hence, term (iii') is zero and by Eq. (S8) term (i) \geq (iii).

Condition (II) directly leads to $I_{X;R|ABS} \geq I_{X;A|BRS}$ as follows.

Define $\mathbf{E} = \{E \in \mathbf{O} \setminus \mathbf{P} : X \not\perp\!\!\!\perp E \mid \mathbf{SO} \setminus \{E\}\}$. By condition (II) there exists a link $E \leftrightarrow W$ or an extended collider path $E \ast \rightarrow C \leftrightarrow \dots \leftrightarrow W$ inside \mathbf{C} for $W \in \mathbf{YM}$ where all colliders $C \in \mathbf{vancs}$. There are two types: (1) $E \rightarrow C \leftrightarrow \dots \leftrightarrow W$ (then $E \in \mathbf{P}_C$) and (2) $E \leftrightarrow W$ or $E \leftrightarrow C \leftrightarrow \dots \leftrightarrow W$. We consider two cases:

Case (1): $E \in \mathbf{E}$ for which there exists *at least one* path of type (1). Any valid \mathbf{Z}' with $E \notin \mathbf{Z}'$ has to block paths from X to E since otherwise there is a non-causal open path from X to Y through the motif chain $\ast \ast E \rightarrow C \leftrightarrow \dots \leftrightarrow W$ for $W \in \mathbf{YM}$: E is open since $E \notin \mathbf{Z}'$ and the part from E to W is open since all colliders $C \in \mathbf{vancs}$: if $C \in an(\mathbf{S})$, the collider is always opened and if $C \in an(XY)$ then either the directed path to X or Y is open, or C is opened if \mathbf{Z}' contains a node on that path.

Case (2): $E \in \mathbf{E}$ for which *all paths* are of type (2). Firstly, all paths from X to E that end with a tail at E must be blocked by \mathbf{Z}' since otherwise there is a non-causal path as for case (1). The same holds for paths that end with a head at E if $E \in \mathbf{vancs}$. Consider paths that end with a head at E and $E \notin \mathbf{vancs}$ which implies $E \perp\!\!\!\perp X \mid \mathbf{vancs}$ by Def. 3. Then it follows that $E \perp\!\!\!\perp X \mid \mathbf{SO} \setminus \{E\}$ and, hence, $E \notin \mathbf{E}$ which we can show by considering where E can occur with respect to the different motifs listed in Lemma B.1 as follows (see the definitions of W, V, F, N, C, P_C there): Motif (1a) “ $\ast \ast X \rightarrow \boxed{C} \leftrightarrow$ ” is not relevant since then non-identifiability holds and motif (2a) “ $X, V \ast \ast \boxed{P} \rightarrow W$ ” is not relevant since $\mathbf{E} \notin \mathbf{P}$. Motifs (3a) “ $V \leftarrow W$ ”, (4a) “ $\ast \ast F \leftrightarrow W$ ”, and (5a) “ $\ast \ast N \leftrightarrow W$ ” are not relevant since no $E \in \mathbf{O}$ is involved. For the motifs (1b) “ $\ast \ast X \rightarrow \boxed{P_C} \rightarrow E \leftrightarrow$ ”, (2b) “ $X, V \ast \ast \boxed{P_C} \rightarrow \boxed{C} \leftrightarrow \dots \leftrightarrow E$ ”, and (3b) “ $X, V \leftarrow \boxed{C} \leftrightarrow \dots \leftrightarrow E$ ” the path to X is blocked by $\mathbf{SO} \setminus \{E\}$. For motif (4b) “ $\ast \ast F \leftrightarrow \boxed{C} \leftrightarrow \dots \leftrightarrow E$ ” or “ $\ast \ast F \leftrightarrow E$ ”, since $\mathbf{SO} \cap \mathbf{forb} = \emptyset$ and $X \cap des(\mathbf{forb}) = \emptyset$, there must exist a collider $\in \mathbf{forb}$ or $\ast \rightarrow F \leftrightarrow$ on a path to X which is then blocked. Hence, $E \notin \mathbf{E}$. Finally, for (5b) “ $\ast \ast N \leftrightarrow \boxed{C} \leftrightarrow \dots \leftrightarrow E$ ” or “ $\ast \ast N \leftrightarrow E$ ” with $N \notin \mathbf{vancs}$ and $X \not\perp\!\!\!\perp N \mid \mathbf{vancs}$ we either have $\ast \rightarrow N \leftrightarrow$ or there exists a collider on any path to X with $K \in des(N)$ and, hence, $K \notin \mathbf{vancs}$. $E \not\perp\!\!\!\perp X \mid \mathbf{SO} \setminus \{E\}$ would only be possible if N or $K \in an(\mathbf{O} \setminus \mathbf{vancs})$. The subset $\mathbf{O} \setminus \mathbf{vancs}$ fulfills $\mathbf{O} \setminus \mathbf{vancs} \perp\!\!\!\perp X \mid \mathbf{vancs}$ by Def. 3. However, N or $K \in an(\mathbf{O} \setminus \mathbf{vancs})$ implies that there is a path from $\mathbf{O} \setminus \mathbf{vancs}$ to N . Then $X \not\perp\!\!\!\perp N \mid \mathbf{vancs}$ contradicts $\mathbf{O} \setminus \mathbf{vancs} \perp\!\!\!\perp X \mid \mathbf{vancs}$ implying that $N, K \notin an(\mathbf{O} \setminus \mathbf{vancs})$ and, hence, $E \perp\!\!\!\perp X \mid \mathbf{SO} \setminus \{E\}$ and $E \notin \mathbf{E}$.

Both cases taken together, it holds that $X \perp\!\!\!\perp E \mid \mathbf{SZ}' \setminus \{E\}$ for any valid \mathbf{Z}' . Furthermore, $X \perp\!\!\!\perp P \mid \mathbf{SZ}' \setminus \{P\}$ with $P \in \mathbf{P}$ for any valid \mathbf{Z}' since \mathbf{P} is directly connected to Y and, therefore, a valid \mathbf{Z}' has to block a non-causal path between X and Y through \mathbf{P} .

Now decompose term (iv) as

$$I_{X;A|Z'S} = \underbrace{I_{X;A_P A_E|Z'S}}_{=0} + I_{X;A \setminus (A_P A_E)|Z'S A_P A_E} \quad (\text{S9})$$

with $A_P = A \cap P$ and $A_E = (A \cap E) \setminus A_P$. The preceding derivations imply $X \perp\!\!\!\perp A_P A_E \mid Z'S$ for any valid \mathbf{Z}' and, hence, the first term vanishes.

Consider the set $\mathbf{E}' = \{E' \in \mathbf{O} \setminus \mathbf{P} : X \perp\!\!\!\perp E' \mid \mathbf{SO} \setminus \{E'\}\}$. This implies that $A_{E'} = A \setminus (A_P A_E)$ fulfills $A_{E'} \perp\!\!\!\perp X \mid \mathbf{SO} \setminus A_{E'}$ and since $\mathbf{SO} \setminus A_{E'} = \mathbf{SBA}_P A_E$ we have

$$I_{X;A_{E'}|\mathbf{SBA}_P A_E} = 0. \quad (\text{S10})$$

This now leads to term (ii) \geq term (iv) by considering two ways of decomposing the following CMI:

$$I_{X;R A_{E'}|\mathbf{SBA}_P A_E} = \underbrace{I_{X;A_{E'}|\mathbf{SBA}_P A_E}}_{=0 \text{ by Eq. (S10)}} + \underbrace{I_{X;R|\mathbf{SBA}_P A_E A_{E'}}}_{\text{term (ii)}} \quad (\text{S11})$$

$$= \underbrace{I_{X;R|\mathbf{SBA}_P A_E}}_{\geq 0} + \underbrace{I_{X;A_{E'}|\mathbf{SBA}_P A_E R}}_{\text{term (iv) by Eq. (S9)}}. \quad (\text{S12})$$

“only if”: We need to prove that if either Condition (I) or Condition (II) or both are not fulfilled, then graphical optimality does not hold (implying that also \mathbf{O} is not optimal).

The negation of Condition (I) directly leads to $I_{\mathbf{A};Y|\mathbf{B}\mathbf{R}\mathbf{X}\mathbf{S}} < I_{\mathbf{R};Y|\mathbf{A}\mathbf{B}\mathbf{X}\mathbf{S}}$ for some distribution \mathcal{P} consistent with \mathcal{G} as follows: There exists at least one N-node with at least one collider path $N \leftrightarrow \dots \leftrightarrow C \leftrightarrow \dots \leftrightarrow W$ (including $N \leftrightarrow W$) for $C \in \mathbf{C}$ and $W \in \mathbf{YM}$ (indexed by i) with collider path nodes denoted π_i^N such that $\mathbf{O}_{\pi_i^N}$ blocks all non-causal paths from X to Y . $\mathbf{O}_{\pi_i^N}$ is the O-set for the causal effect of X on Y given $\mathbf{S}' = \mathbf{S} \cup \{N\} \cup \pi_i^N$. Consider $\mathbf{Z}' = \mathbf{O}_{\pi_i^N}$. Since also $N \notin \mathbf{forb}$, $\mathbf{S} \cap \mathbf{forb} = \emptyset$, and $\pi_i^N \cap \mathbf{forb} = \emptyset$, $\mathbf{Z}' = \mathbf{O}_{\pi_i^N}$ is valid. Since $N \in \mathbf{O}_{\pi_i^N}$ while $N \notin \mathbf{O}$, we have $\mathbf{R} \neq \emptyset$, and since $\pi_i^N \subseteq \mathbf{C}$ we have $\pi_i^N \subseteq \mathbf{B}\mathbf{S}$ and there exists an open path $N \leftrightarrow \boxed{C} \leftrightarrow \dots \leftrightarrow \boxed{C} \leftrightarrow W$ (or $N \leftrightarrow W$) such that $I_{\mathbf{R};Y|\mathbf{B}\mathbf{X}\mathbf{S}} > 0$. Similar to Lemma 1 we can now construct a distribution \mathcal{P} with associated SCM consistent with \mathcal{G} where all links “ $U \ast \ast A$ ” for $A \in \mathbf{A}$ almost vanish and, hence, term (i’), $I_{\mathbf{A};Y|\mathbf{B}\mathbf{X}\mathbf{S}} \rightarrow 0$: Consider the three possible links and associated arbitrary assignment functions in the SCM: (1) “ $A \rightarrow U$ ” with $U := f_U(\dots, cA, \dots)$, (2) “ $A \leftarrow U$ ” with $A := f_A(\dots, cU, \dots)$, and (3) “ $A \leftrightarrow U$ ” with $A := f_A(\dots, cL^U, \dots)$ where L^U denotes one or more latent variables and $c \in \mathbb{R}$. We then consider the limit $c \rightarrow 0$ leading to term (i’), $I_{\mathbf{A};Y|\mathbf{B}\mathbf{X}\mathbf{S}} \rightarrow 0$. Since $\mathbf{A} \cap N\pi_i^N = \emptyset$ this does not affect the collider path $N \leftrightarrow \boxed{C} \leftrightarrow \dots \leftrightarrow \boxed{C} \leftrightarrow W$ (or $N \leftrightarrow W$) such that $I_{\mathbf{R};Y|\mathbf{B}\mathbf{X}\mathbf{S}} > 0$. By Eq. (S8) then term (i)<(iii). By Lemma 1, where \mathcal{P} is further modified to \mathcal{P}' without affecting term (i)<(iii), then graphical optimality does not hold.

Alternatively, the negation of Condition (II) directly leads to $I_{X;\mathbf{R}|\mathbf{A}\mathbf{B}\mathbf{S}} < I_{X;\mathbf{A}|\mathbf{B}\mathbf{R}\mathbf{S}}$ as follows: By the negation of Condition (II) there exists an $E \in \mathbf{O} \setminus \mathbf{P}$ with $X \not\leftrightarrow E | \mathbf{S}\mathbf{O} \setminus \{E\}$ such that there is no link $E \leftrightarrow W$ and all extended collider paths $E \ast \ast C \leftrightarrow \dots \leftrightarrow W$ inside \mathbf{C} for $W \in \mathbf{YM}$ contain at least one collider $C \notin \mathbf{vancs}$. Define the set of these non-ancestral colliders as

$$\mathbf{C}_E = \{C \in \mathbf{C} : E \ast \ast \dots \leftrightarrow C \leftrightarrow \dots \leftrightarrow W\} \setminus \mathbf{vancs}. \quad (\text{S13})$$

We define $E_C = \{E\} \cup (\text{des}(\mathbf{C}_E) \cap \mathbf{O})$ and choose $\mathbf{Z}' = \mathbf{O} \setminus E_C$ implying $\mathbf{A} = E_C$, $\mathbf{B} = \mathbf{O} \setminus E_C$, and $\mathbf{R} = \emptyset$. We need to show that (1) \mathbf{Z}' is valid and (2) $I_{X;\mathbf{A}|\mathbf{B}\mathbf{R}\mathbf{S}} = I_{X;E_C|\mathbf{S}\mathbf{O} \setminus E_C} > I_{X;\mathbf{R}|\mathbf{A}\mathbf{B}\mathbf{S}} = 0$ (since $\mathbf{R} = \emptyset$).

Ad (1): As a subset of \mathbf{O} we have that $\mathbf{Z}' \cap \mathbf{forb} = \emptyset$. We investigate whether \mathbf{Z}' blocks all non-causal paths between X and Y by considering the motifs in Lemma B.1. In addition to all those motifs listed there there are modified motifs where unconditioned C -nodes and P_C -nodes occur (denoted without a $\boxed{\cdot}$) due to removing E_C from \mathbf{O} .

Firstly, the unmodified motifs are blocked as before (see Theorem 1): Motif (1a) “ $\ast \ast X \rightarrow \boxed{C} \leftrightarrow$ ” is not relevant since then non-identifiability holds. By Lemma B.4 the motifs (1b), (2a), (2b), (3a), (3b) all contain a non-collider in $\mathbf{S}\mathbf{O} \setminus E_C$ and are blocked. By Lemma B.5 (part one) the motifs (4a,b) for $F \in \text{des}(\mathbf{YM})$ are blocked because $\mathbf{Z}' \cap \mathbf{forb} = \emptyset$. By Lemma B.6 (part one) motifs (5a) and (5b) are blocked given $\mathbf{S}\mathbf{O} \setminus E_C$ because the proof in Lemma B.6 requires that on paths to X either N is a collider or there exists a descendant collider K and that $N, K \notin \text{an}(\mathbf{O}\mathbf{S})$. The latter is fulfilled because $\mathbf{S}\mathbf{O} \setminus E_C$ is a subset of $\mathbf{S}\mathbf{O}$.

Secondly, all paths from X through the removed node E to $W \in \mathbf{YM}$ are blocked by $\mathbf{S}\mathbf{O} \setminus E_C$: Paths through \mathbf{P} are blocked since $E \notin \mathbf{P}$ and $\text{des}(\mathbf{C}_E) \cap \mathbf{vancs} = \emptyset$ and, hence, $\mathbf{P} \subseteq \mathbf{S}\mathbf{O} \setminus E_C$. Paths through colliders are blocked by the negation of condition (II): there is no link $E \leftrightarrow W$ and all extended collider paths $E \ast \ast C \leftrightarrow \dots \leftrightarrow W$ inside \mathbf{C} for $W \in \mathbf{YM}$ contain at least one collider $C \notin \mathbf{vancs}$. By construction, $E_C = \{E\} \cup (\text{des}(\mathbf{C}_E) \cap \mathbf{O})$, implying that all these non-ancestral colliders are blocked.

Thirdly, we consider the modified motifs with unconditioned $C, P_C \in (\text{des}(\mathbf{C}_E) \cap \mathbf{O})$. By definition of \mathbf{C}_E in (S13), $C, P_C \notin \mathbf{vancs}$. (As a remark, E can potentially be in \mathbf{vancs} .) Motif (1a) “ $\ast \ast X \rightarrow \boxed{C} \leftrightarrow$ ” cannot occur since then non-identifiability holds. Modified motifs (1,2b’) “ $X, V \ast \ast \boxed{P_C} \rightarrow C \leftrightarrow$ ” and (1,2b’’) “ $X, V \ast \ast \boxed{P_C} \rightarrow \boxed{C} \leftrightarrow \dots \leftrightarrow C \leftrightarrow$ ” are blocked since they contain a non-collider. Motifs (1,2b’’’) “ $X, V \ast \ast P_C \rightarrow C \leftrightarrow$ ” are blocked since $C \notin \text{des}(\mathbf{S}\mathbf{O} \setminus E_C)$. Motif (2a’) “ $X, V \ast \ast P \rightarrow W$ ” is not possible since $P \in \mathbf{P} \subseteq \mathbf{vancs}$. Motifs (3a) “ $V \leftarrow W$ ”, (4a) “ $\ast \ast F \leftrightarrow W$ ”, and (5a) “ $\ast \ast N \leftrightarrow W$ ” are not modified since no conditioned node occurs. Motif (3b’) “ $X, V \leftarrow C \leftrightarrow$ ” is blocked because due to $C \notin \mathbf{vancs}$ there must exist a descendant of C that is a collider $K \notin \mathbf{vancs}$ on the path to X . Since E_C contains all descendants of C , also K and all its descendants are not in $\mathbf{S}\mathbf{O} \setminus E_C$ and K is blocked. Finally, motifs (4b’) “ $\ast \ast F \leftrightarrow C \leftrightarrow$ ” and (5b’) “ $\ast \ast N \leftrightarrow C \leftrightarrow$ ” are blocked since $\mathbf{S}\mathbf{O} \setminus E_C$ does not contain any descendant of C . This proves the validity of \mathbf{Z}' .

Ad (2): To show that $I_{X;\mathbf{A}|\mathbf{BRS}} = I_{X;E_C|\mathbf{SO}\setminus E_C} > I_{X;\mathbf{R}|\mathbf{ABS}} = 0$ we start from the assumption that $E \not\perp X \mid \mathbf{SO} \setminus \{E\}$. This implies that there exists a path from X to E where no non-collider is in $\mathbf{SO} \setminus \{E\}$ and for every collider K it holds that $des(K) \cap \mathbf{SO} \setminus \{E\} \neq \emptyset$. With $\mathbf{Z}' = \mathbf{O} \setminus E_C$ all non-colliders are still open. Consider those colliders K with $des(K) \cap (\mathbf{O} \setminus \{E\}) \subseteq E_C \setminus \{E\}$. Then these colliders are closed on the path from E to X . However, for each such K there is a $C \in E_C \setminus \{E\}$ with $C \in des(K)$. Then the path from X through $* \rightarrow K \rightarrow \dots \rightarrow C$ is open given $\mathbf{SO} \setminus E_C$. Hence, at least for the last such collider on the path from E to X there is an open path from $C \in E_C \setminus \{E\}$ to X given $\mathbf{SO} \setminus E_C$. Then Faithfulness implies that $I_{X;\mathbf{A}|\mathbf{BRS}} = I_{X;E_C|\mathbf{SO}\setminus E_C} > I_{X;\mathbf{R}|\mathbf{ABS}} = 0$ and, hence, term (ii) < term (iv) holds for all distributions \mathcal{P} consistent with \mathcal{G} . By Lemma 1, where the distribution \mathcal{P} is modified to \mathcal{P}' without affecting term (ii)<(iv), then graphical optimality does not hold.

This concludes the proof of Theorem 3. \square

B.11 Proof of Corollary B.1

Corollary (Minimality and minimum cardinality). Given Assumptions 1, assume that graphical optimality holds, and, hence, \mathbf{O} is optimal. Further it holds that:

1. If \mathbf{O} is not minimal, then $J_{\mathbf{O}} > J_{\mathbf{Z}}$ for all *minimal* valid $\mathbf{Z} \neq \mathbf{O}$,
2. If \mathbf{O} is minimal valid, then \mathbf{O} is the unique set that maximizes the adjustment information $J_{\mathbf{Z}}$ among all *minimal* valid $\mathbf{Z} \neq \mathbf{O}$,
3. \mathbf{O} is of minimum cardinality, that is, there is no subset of \mathbf{O} that is still valid and optimal.

Proof. We again define disjunct sets $\mathbf{R}, \mathbf{B}, \mathbf{A}$ with $\mathbf{A} = \mathbf{O} \setminus \mathbf{Z}$, $\mathbf{R} = \mathbf{Z} \setminus \mathbf{O}$, and $\mathbf{B} = \mathbf{O} \cap \mathbf{Z}$, where any of them can be empty, but not both \mathbf{R} and \mathbf{A} since then $\mathbf{Z} = \mathbf{O}$. Hence $\mathbf{O} = \mathbf{AB}$ and $\mathbf{Z} = \mathbf{BR}$. Consider relation (S5) in this case,

$$J_{\mathbf{O}} = J_{\mathbf{Z}} + \underbrace{I_{\mathbf{A};Y|\mathbf{BRXS}}}_{(i)} + \underbrace{I_{X;\mathbf{R}|\mathbf{ABS}}}_{(ii)} - \underbrace{I_{\mathbf{R};Y|\mathbf{ABXS}}}_{(iii)} - \underbrace{I_{X;\mathbf{A}|\mathbf{BRS}}}_{(iv)}. \quad (\text{S14})$$

Part 1 and 2: Since graphical optimality holds, we know that $J_{\mathbf{O}} = J_{\mathbf{Z}}$ can only be achieved if term (i) = term (iii) and term (ii) = term (iv). From Eq. (S8) we know that term (i) = (iii) can only hold if $I_{\mathbf{A};Y|\mathbf{BXS}} = 0$. But this implies $\mathbf{A} = \emptyset$ by Faithfulness since, by construction, $\mathbf{A} \subset \mathbf{O}$ is always connected to Y (potentially through \mathbf{M}) given $X\mathbf{SO} \setminus \mathbf{A}$. Then term (iv) = 0 and, by optimality, $I_{X;\mathbf{R}|\emptyset\mathbf{BS}} = 0$. But the latter would imply that $\mathbf{Z} = \mathbf{BR}$ is either not minimal anymore since \mathbf{R} is not connected to X and, hence, does not block any non-causal path not already blocked by \mathbf{B} . Then $J_{\mathbf{O}} > J_{\mathbf{Z}}$ among all minimal valid \mathbf{Z} (Part 1). Or \mathbf{Z} is minimal and $\mathbf{R} = \emptyset$, for which $\mathbf{Z} = \mathbf{O}$ is the unique set maximizing $J_{\mathbf{Z}}$ among all minimal valid $\mathbf{Z} \neq \mathbf{O}$ (Part 2).

Part 3, i.e., that removing any subset from \mathbf{O} decreases $J_{\mathbf{O}}$ follows directly from setting $\mathbf{R} = \emptyset$ and considering $\mathbf{A} \neq \emptyset$ (since otherwise nothing would be removed). Then term (ii) and term (iii) are both zero and by optimality term (iv), which must be smaller or equal to term (ii), is zero. Since \mathbf{A} is connected to Y (see Part 1) by Faithfulness we have $J_{\mathbf{O}} > J_{\mathbf{O}\setminus\mathbf{A}}$. \square

C Algorithms

Algorithm C.1 Construction of \mathbf{O} -set and test for backdoor-identifiability.

Require: Causal graph \mathcal{G} , cause variable X , effect variable Y , mediators \mathbf{M} , conditioned variables \mathbf{S}

- 1: Initialize $\mathbf{P} = \emptyset$, $\mathbf{C} = \emptyset$ and $\mathbf{P}_{\mathbf{C}} = \emptyset$
- 2: **for** $W \in \mathbf{YM}$ **do**
- 3: $\mathbf{P} = \mathbf{P} \cup pa(W) \setminus \text{forb}$
- 4: **for** $W \in \mathbf{YM}$ **do**
- 5: Initialize nodes in this level $\mathcal{L} = \{W\}$
- 6: Initialize ignorable nodes $\mathcal{N} = \emptyset$
- 7: **while** $|\mathcal{L}| > 0$ **do**
- 8: Initialize next level $\mathcal{L}' = \emptyset$
- 9: **for** $C \in sp(\mathcal{L}) \setminus \mathcal{N}$ **do**
- 10: **if** $C = X$ **then**
- 11: **return** No valid backdoor adjustment set exists.
- 12: **if** $C \notin \mathbf{C}$ and Def. 3 (1) $C \notin \text{forb}$ and ((2a) $C \in \text{vancs}$ or (2b) $C \perp\!\!\!\perp X \mid \text{vancs}$)
- 13: **then**
- 14: $\mathbf{C} = \mathbf{C} \cup \{C\}$
- 15: $\mathcal{L}' = \mathcal{L}' \cup \{C\}$
- 16: **else**
- 17: **if** $C \notin \mathbf{C}$ **then**
- 18: $\mathcal{N} = \mathcal{N} \cup \{C\}$
- 19: $\mathcal{L} = \mathcal{L}' \setminus \mathcal{N}$
- 20: **for** $C \in \mathbf{C}$ **do**
- 21: **if** $X \in pa(C)$ **then**
- 22: **return** No valid backdoor adjustment set exists.
- 23: $\mathbf{P}_{\mathbf{C}} = \mathbf{P}_{\mathbf{C}} \cup pa(C)$
- 24: **return** $\mathbf{O} = \mathbf{PCP}_{\mathbf{C}}$

Algorithm C.2 Construction of \mathbf{O}_{\min} and $\mathbf{O}_{\mathbf{C}_{\min}}$ -sets. The relevant code for $\mathbf{O}_{\mathbf{C}_{\min}}$ is indicated in parentheses.

Require: Causal graph \mathcal{G} , cause variable X , effect variable Y , mediators \mathbf{M} , conditioned variables \mathbf{S} , $\mathbf{O} = \mathbf{PCP}_{\mathbf{C}}$ -set

- 1: Initialize $\mathbf{O}_{\min} = \mathbf{O}$ ($\mathbf{C}_{\min} = \mathbf{C}_{\mathbf{P}} \setminus \mathbf{P}$)
- 2: **for** $Z \in \mathbf{O}_{\min}$ ($Z \in \mathbf{C}_{\min}$) **do**
- 3: **if** Z has no active path to X given $\mathbf{SO} \setminus \{Z\}$ **then**
- 4: Mark Z for removal
- 5: Remove marked nodes from \mathbf{O}_{\min} (\mathbf{C}_{\min})
- 6: **for** $Z \in \mathbf{O}_{\min}$ ($Z \in \mathbf{C}_{\min}$) **do**
- 7: **if** Z has no active path to Y given $X\mathbf{SO}_{\min} \setminus \{Z\}$ (given $X\mathbf{SPC}_{\min} \setminus \{Z\}$) **then**
- 8: Mark Z for removal
- 9: Remove marked nodes from \mathbf{O}_{\min} (\mathbf{C}_{\min})
- 10: **return** \mathbf{O}_{\min} ($\mathbf{O}_{\mathbf{C}_{\min}} = \mathbf{PC}_{\min}$)

D Further details and figures of further numerical experiments

D.1 Setup

We compare the following adjustment sets (see definitions in Section 2.3):

- \mathbf{O}
- Adjust
- $\mathbf{O}_{C_{\min}}$
- \mathbf{O}_{\min}
- $\text{Adjust}_{X_{\min}}$
- Adjust_{\min}

To investigate the applicability of different estimators, we use above adjustment sets together with the following estimators from `sklearn` (version 0.24.2) and the `doublem1` (version 0.4.0) package (see instantiated class for parameters):

- Linear ordinary least squares (LinReg) regressor `LinearRegression()`
- k -nearest-neighbor (kNN) regressor `KNeighborsRegressor(n_neighbors=3)`
- Multilayer perceptron (MLP) regressor `MLPRegressor(max_iter=2000)`
- Random forest (RF) [Breiman, 2001] regressor `RandomForestRegressor()`
- Double machine learning for partially linear regression models (DML) [Chernozhukov et al., 2018] `DoubleMLPLR(data, ml_g, ml_m)` from `doublem1` with `ml_g=ml_g=MLPRegressor(max_iter=2000)` from `sklearn`

`Sklearn` [Pedregosa et al., 2011] and `doublem1` [Bach et al., 2021] are both available under an MIT license.

As data generating processes we consider linear and nonlinear experiments generated with the following generalized additive model:

$$V^j = \sum_i c_i f_i(V^i) + \eta^j \quad \text{for } j \in \{1, \dots, \tilde{N}\}. \quad (\text{S15})$$

To generate a structural causal model among \tilde{N} variables we randomly choose L links whose functional dependencies are linear for linear experiments and one half is $f_i(x) = (1 + 5xe^{-x^2/20})x$ for nonlinear experiments. Coefficients c_i are drawn uniformly from $\pm[0.1, 2]$. For linear experiments we use normal noise $\eta^j \sim \mathcal{N}(0, \sigma^2)$ and, in addition, for nonlinear models $\frac{1}{3}$ of the noise terms is Weibull-distributed, both with standard deviation σ drawn uniformly from $[0.5, 2]$. From the \tilde{N} variables of each dataset we randomly choose a fraction λ as unobserved and denote the number of observed variables as N . For each combination of $N \in \{5, 10, 15, 20\}$, $L \in \{2\tilde{N}, 3\tilde{N}\}$, and $\lambda \in \{30\%, 40\%, 50\%\}$ we randomly create a structural causal model and then randomly pick an observed pair $(X = V^i, Y = V^j)$ connected by a causal path, set $\mathbf{S} = \emptyset$, and consider the intervention $do(V^i = V^i + 1 = x)$ relative to the unperturbed data (x') as ground truth, which corresponds to the linear regression coefficient in the linear case. We further assert that the following criteria hold: (1) the effect is identifiable, (2) the minimal adjustment cardinality is $|\text{vancs}_{\min}(X, Y)| > 0$, and (3) the (absolute) causal effect is $\geq 10^{-3}$ to make sure that Faithfulness holds (if these criteria cannot be fulfilled, another model is generated). We create 500 models for each combination of N, L, λ . Surprisingly, among in total 12,000 randomly created configurations 93% fulfill the optimality conditions in Thm. 3. This may indicate that also in many real-world scenarios graphical optimality actually holds. Here we do not consider the effect of a selected conditioning variable \mathbf{S} since it would have a similar effect on all methods considered.

For the considered graphs the computation time to construct adjustment sets is very short and arguably negligible to the actual cost of fitting methods that use these adjustment sets. The results were evaluated on Intel Xeon Platinum 8260.

D.2 Figures for linear least squares estimator

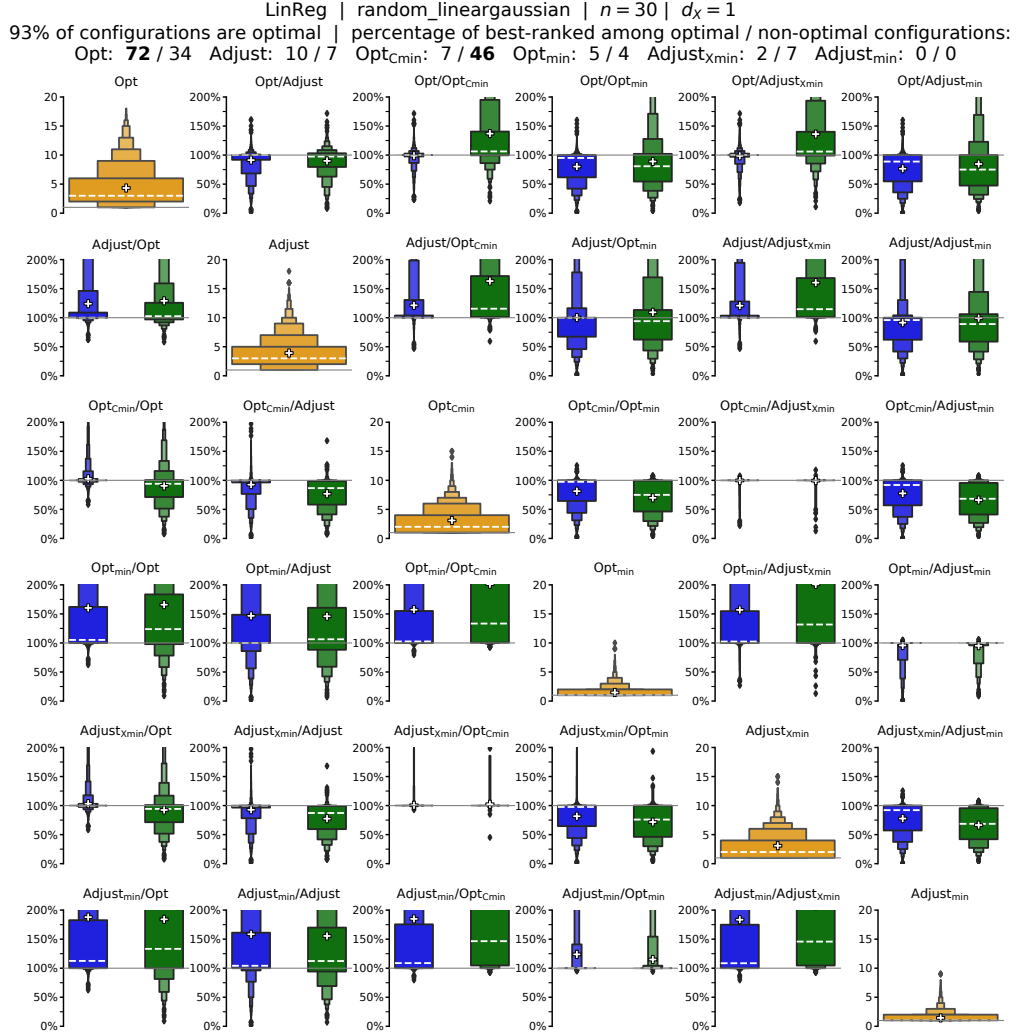


Figure S2: Results of linear experiments with linear estimator and sample size $n = 30$. The diagonal depicts letter-value plots [Hofmann et al., 2017] of adjustment set cardinalities and the off-diagonal shows pairs of RMSE ratios for all combinations of (O, Adjust, O_{Cmin}, O_{min}, Adjust_{Xmin}, Adjust_{min}) for optimal configurations (left in blue) and non-optimal configurations (right in green). Values above 200% are not shown. The dashed horizontal line denotes the median of the RMSE ratios, and the white plus their average. The letter-value plots are interpreted as follows: The largest box shows the 25%–75% range. The next smaller box above (below) shows the 75%–87.5% (12.5%–25%) range and so forth. The numbers on best-ranked methods at the top indicate the percentage of the 12,000 randomly created configurations where the method had the lowest variance. The highest percentage is marked in bold. Note that the highest ranked method may outperform others only by a small margin. The results in the letter-value plots provide a more quantitative picture. See also Fig. S7 where the ranks are further distinguished by the O-set cardinality.

LinReg | random_lineargaussian | $n = 50$ | $d_x = 1$
 93% of configurations are optimal | percentage of best-ranked among optimal / non-optimal configurations:
 Opt: 78 / 39 Adjust: 10 / 6 Opt_{Cmin}: 5 / 45 Opt_{min}: 3 / 3 Adjust_{xmin}: 2 / 5 Adjust_{min}: 0 / 0

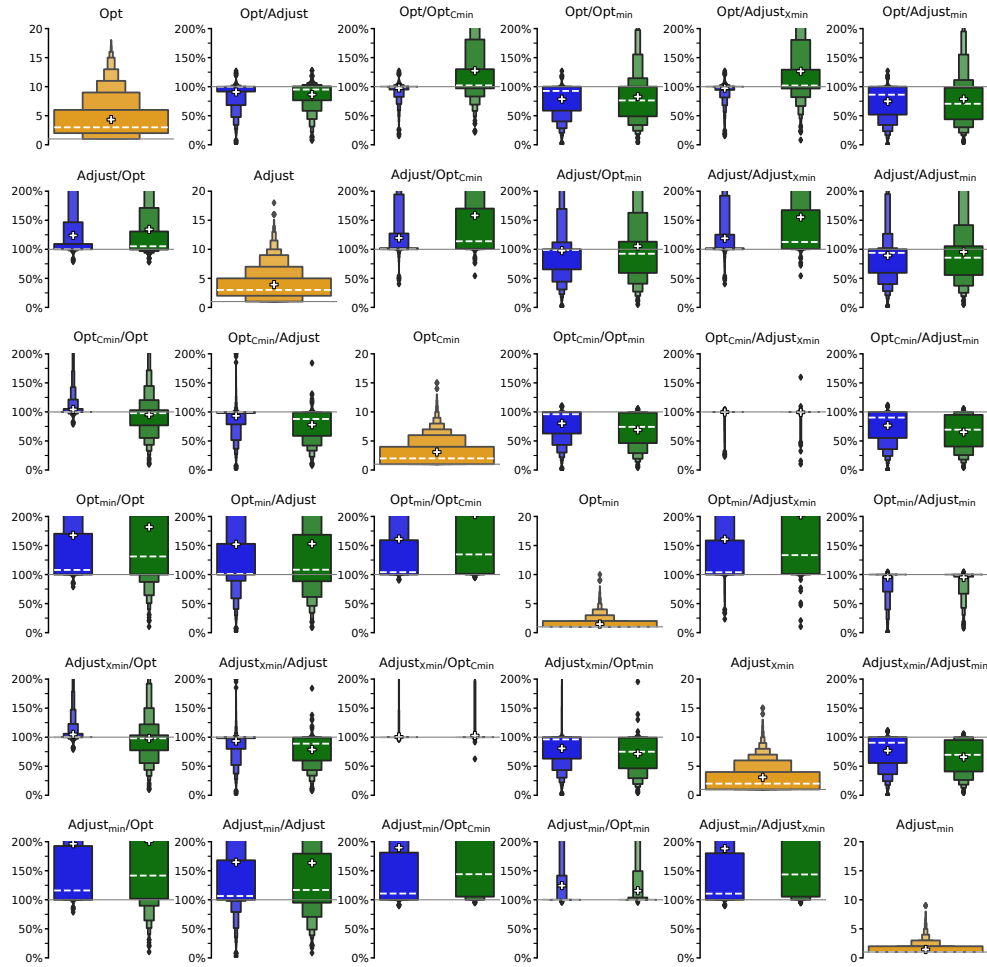


Figure S3: As in Fig. S2 but for $n = 50$.

LinReg | random_lineargaussian | $n = 100$ | $d_x = 1$
 93% of configurations are optimal | percentage of best-ranked among optimal / non-optimal configurations:
 Opt: **81 / 45** Adjust: 9 / 6 Opt_{Cmin}: 4 / 38 Opt_{min}: 2 / 2 Adjust_{xmin}: 1 / 6 Adjust_{min}: 0 / 0

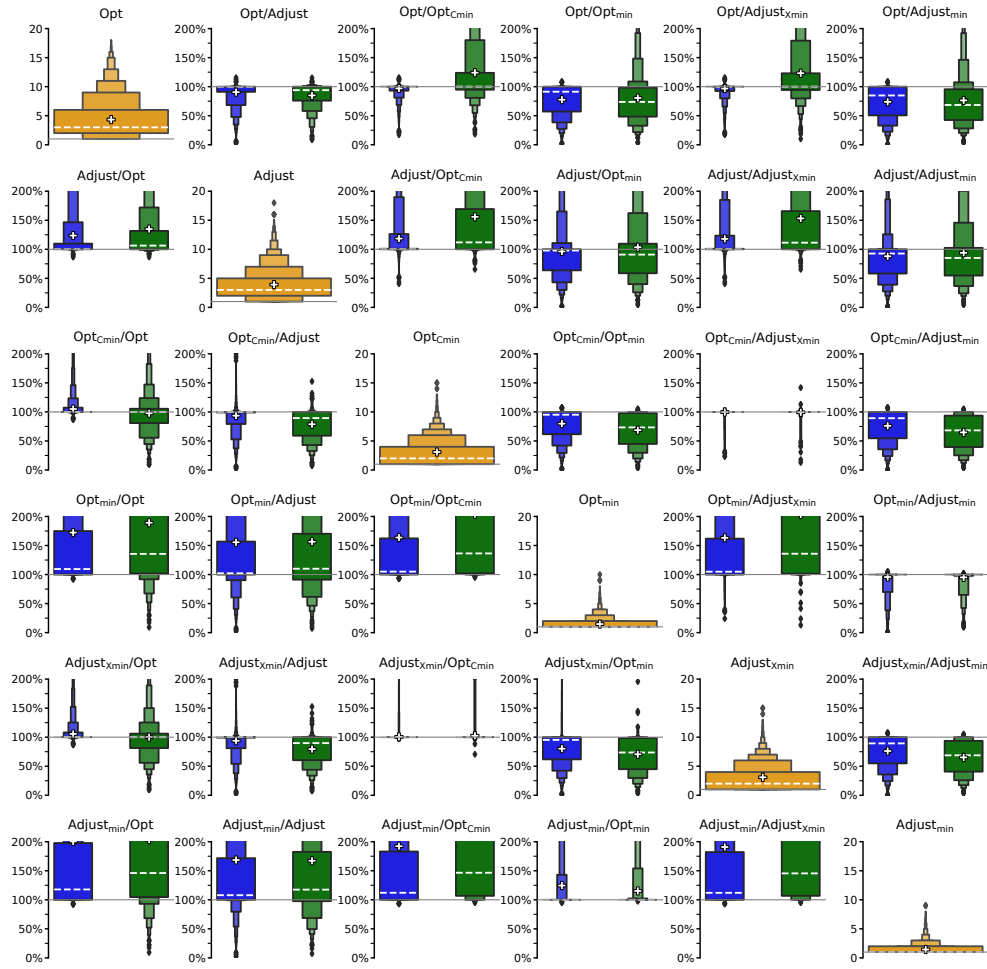


Figure S4: As in Fig. S2 but for $n = 100$.

LinReg | random_lineargaussian | $n = 1000$ | $d_X = 1$
 93% of configurations are optimal | percentage of best-ranked among optimal / non-optimal configurations:
 Opt: **85 / 50** Adjust: 8 / 6 Opt_{Cmin}: 2 / 36 Opt_{min}: 1 / 0 Adjust_{xmin}: 1 / 5 Adjust_{min}: 0 / 0

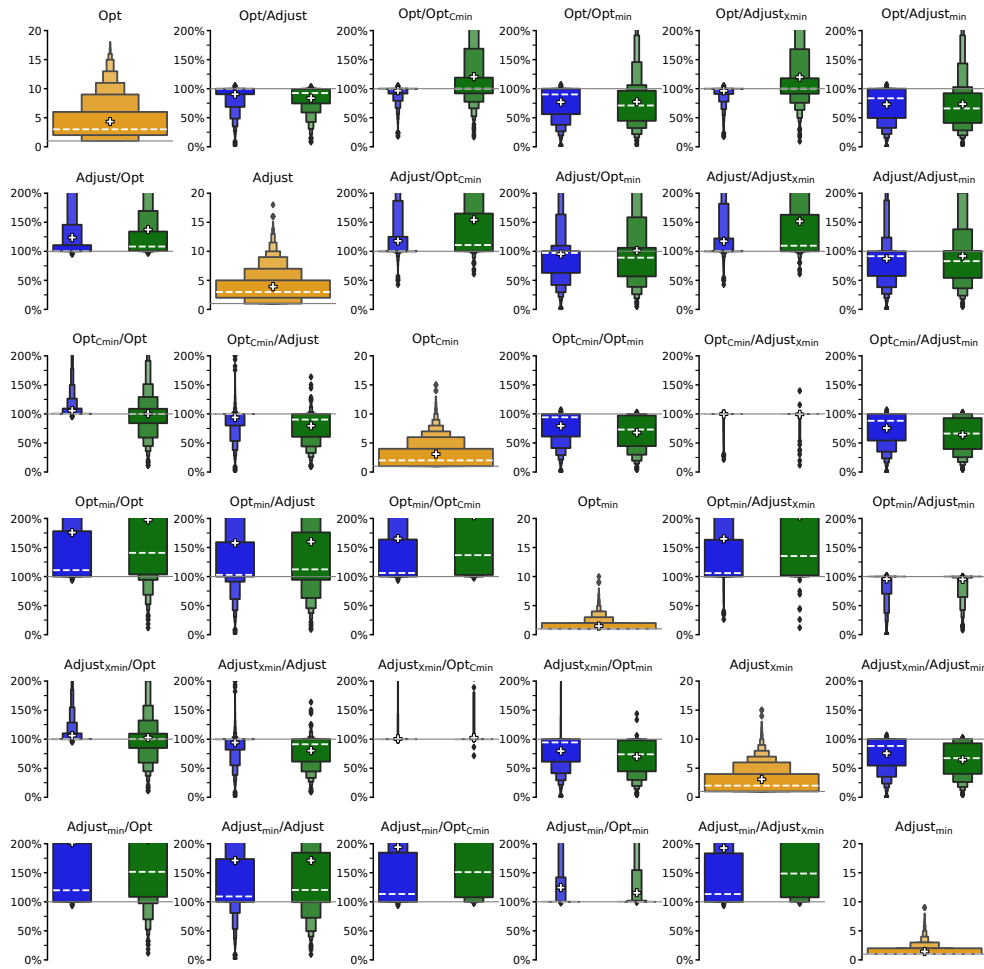


Figure S5: As in Fig. S2 but for $n = 1000$.

LinReg | random_lineargaussian | $n = 10000$ | $d_X = 1$
 93% of configurations are optimal | percentage of best-ranked among optimal / non-optimal configurations:
 Opt: **86 / 53** Adjust: 8 / 4 Opt_{Cmin}: 1 / 36 Opt_{min}: 1 / 0 Adjust_{xmin}: 1 / 5 Adjust_{min}: 0 / 0

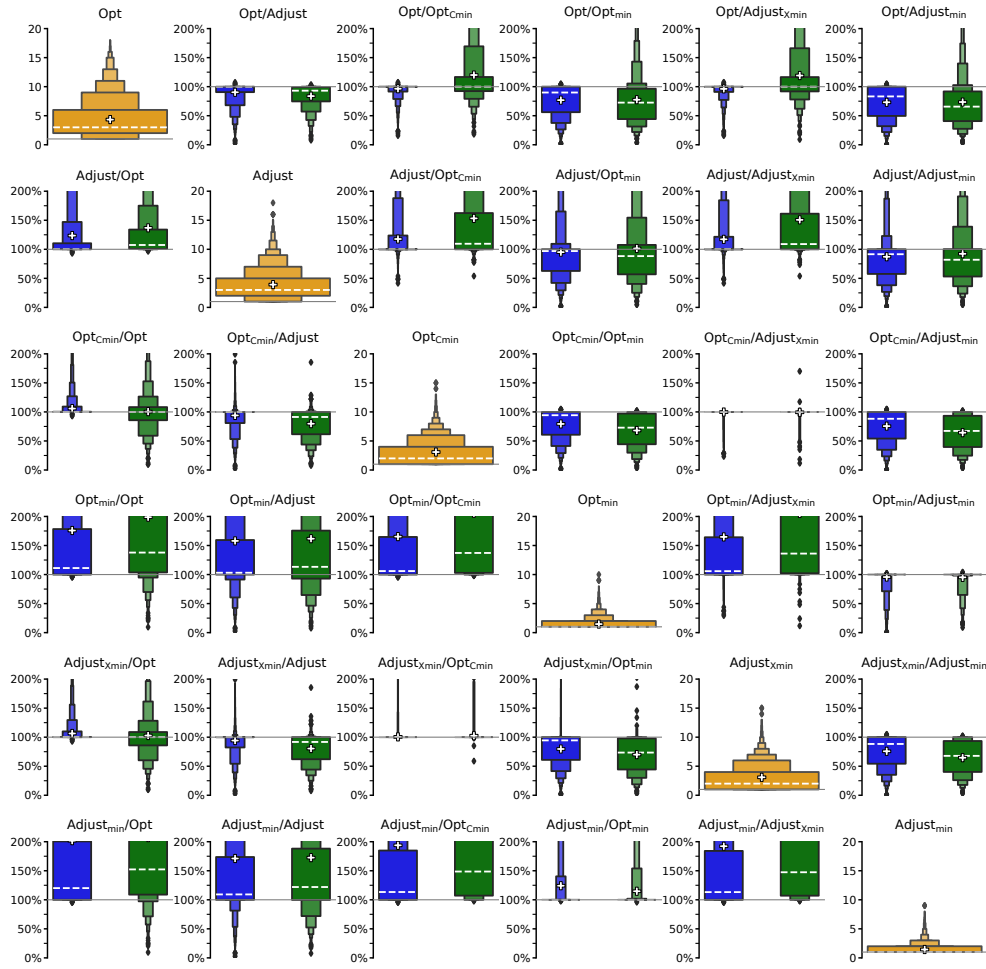


Figure S6: As in Fig. S2 but for $n = 10000$.

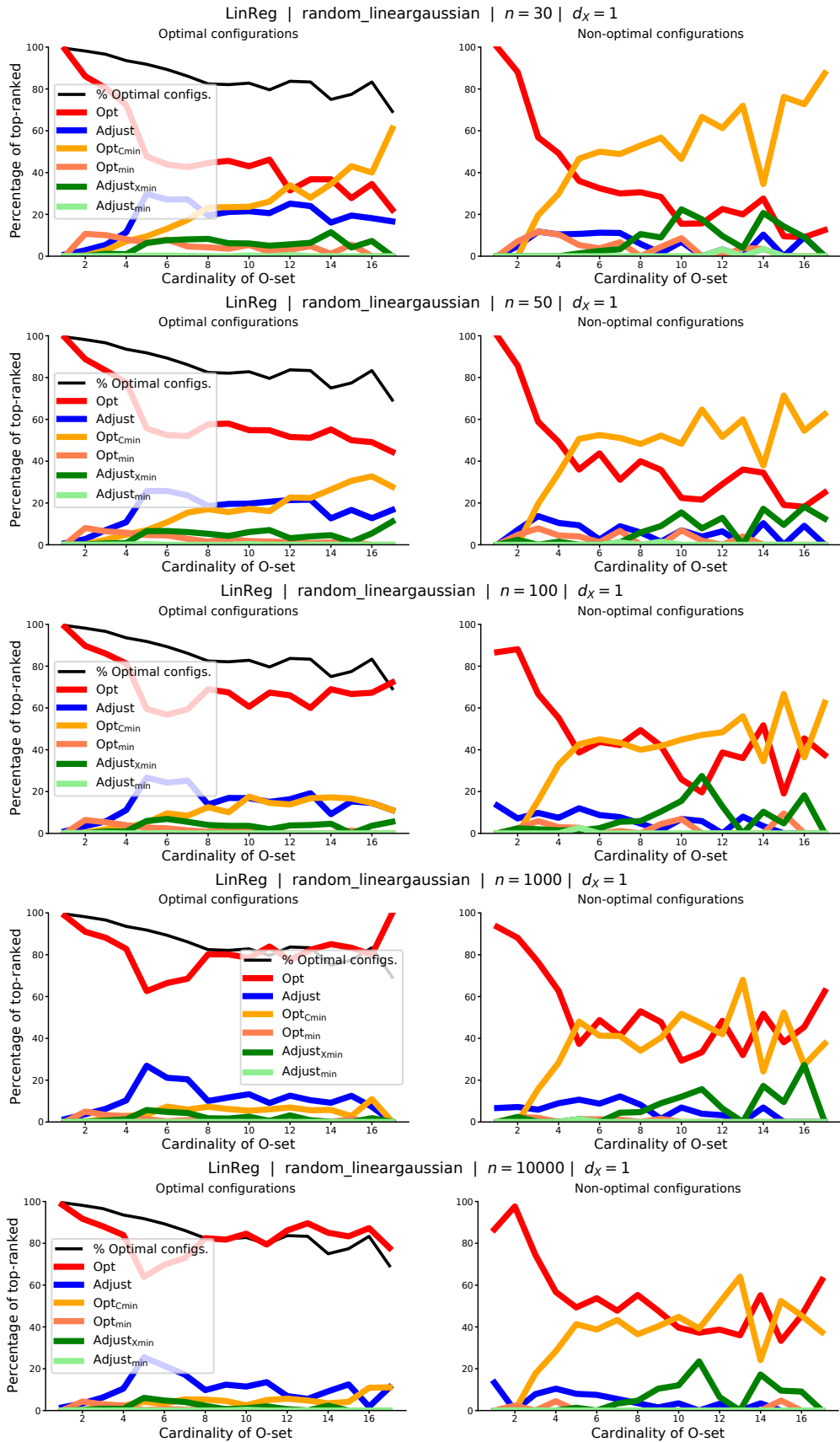


Figure S7: Percentage of configurations where each method has the lowest variance for linear experiments, stratified by the cardinality of the O -set (x -axis) for $n = 30$ (top) to $n = 10,000$ (bottom).

D.3 Figures for non-parametric estimators

kNN | random_lineargaussian | $n = 1000$ | $d_X = 1$
 93% of configurations are optimal | percentage of best-ranked among optimal / non-optimal configurations:
 Opt: 46 / 11 Adjust: 7 / 8 Opt_{Cmin}: 4 / 25 Opt_{min}: 39 / 50 Adjust_{Xmin}: 0 / 0 Adjust_{min}: 2 / 4

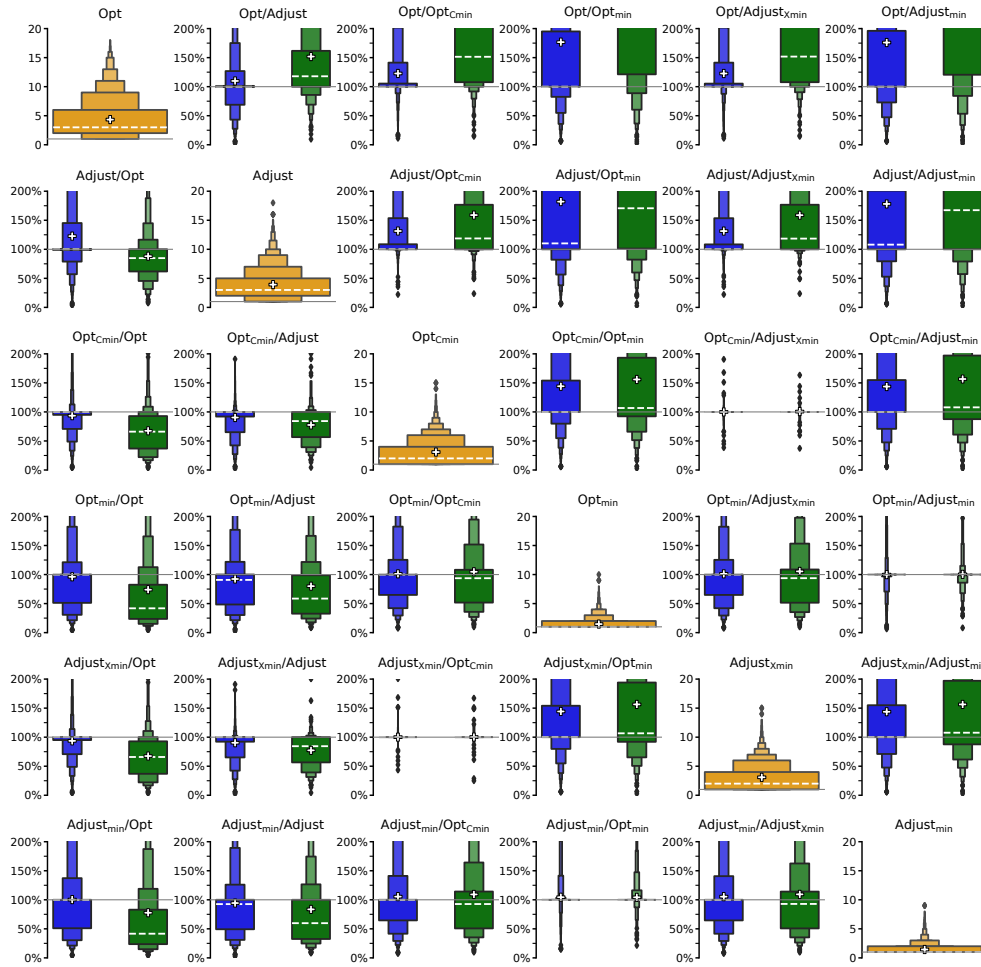


Figure S8: As in Fig. S2 but with kNN estimator ($k = 3$) and $n = 1000$. See also Figs. S15,S16 where the ranks are further distinguished by the O-set cardinality.

kNN | random_nonlinearmixed | $n = 1000$ | $d_X = 1$
 93% of configurations are optimal | percentage of best-ranked among optimal / non-optimal configurations:
 Opt: 55 / 22 Adjust: 11 / 14 Opt_{Cmin}: 5 / 31 Opt_{min}: 25 / 28 Adjust_{xmin}: 0 / 0 Adjust_{min}: 2 / 2

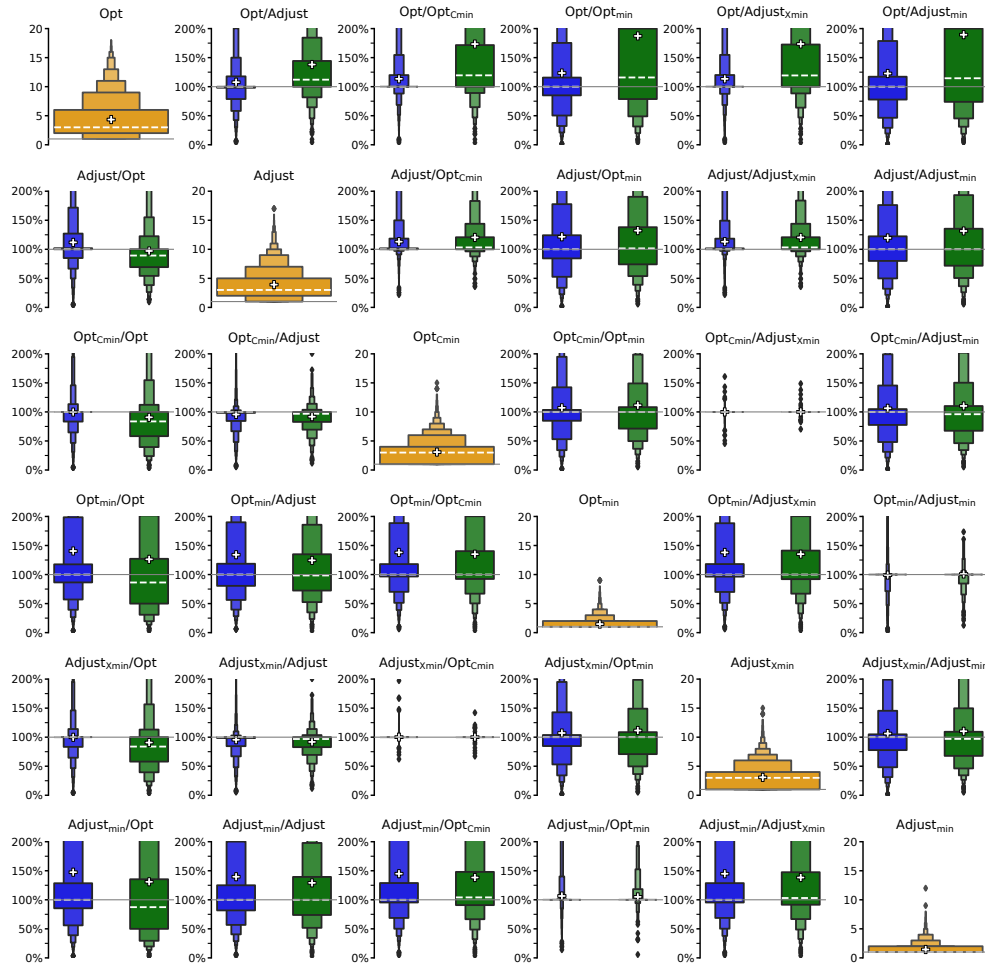


Figure S9: As in Fig. S2 but for kNN estimator ($k = 3$), the nonlinear model, and $n = 1000$.

MLP | random_lineargaussian | $n = 1000$ | $d_X = 1$
 93% of configurations are optimal | percentage of best-ranked among optimal / non-optimal configurations:
 Opt: 31 / 32 Adjust: 13 / 4 Opt_{Cmin}: 18 / 26 Opt_{min}: 10 / 6 Adjust_{Xmin}: 18 / 25 Adjust_{min}: 8 / 5

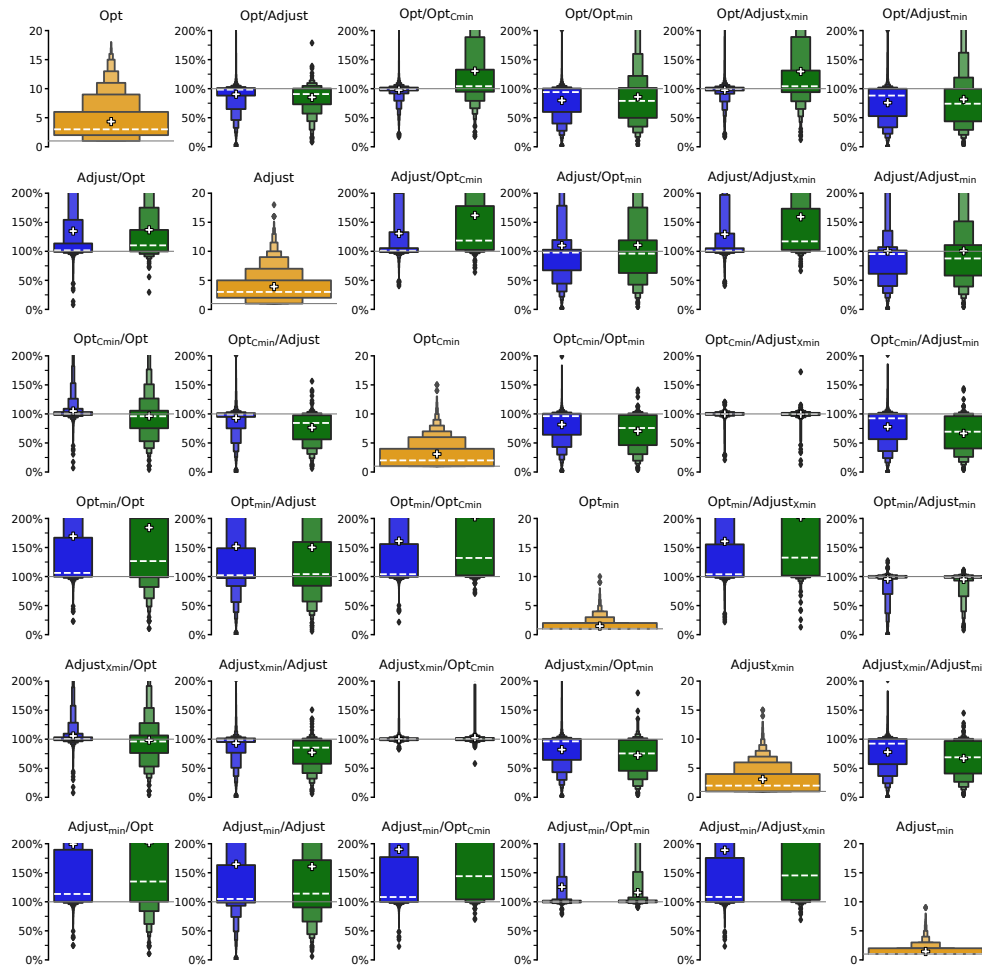


Figure S10: As in Fig. S2 but for MLP estimator and $n = 1000$.

MLP | random_nonlinearmixed | $n = 1000$ | $d_X = 1$
 93% of configurations are optimal | percentage of best-ranked among optimal / non-optimal configurations:
 Opt: 24 / 28 Adjust: 15 / 9 Opt_{Cmin}: 17 / 24 Opt_{min}: 13 / 9 Adjust_{xmin}: 17 / 19 Adjust_{min}: 11 / 9

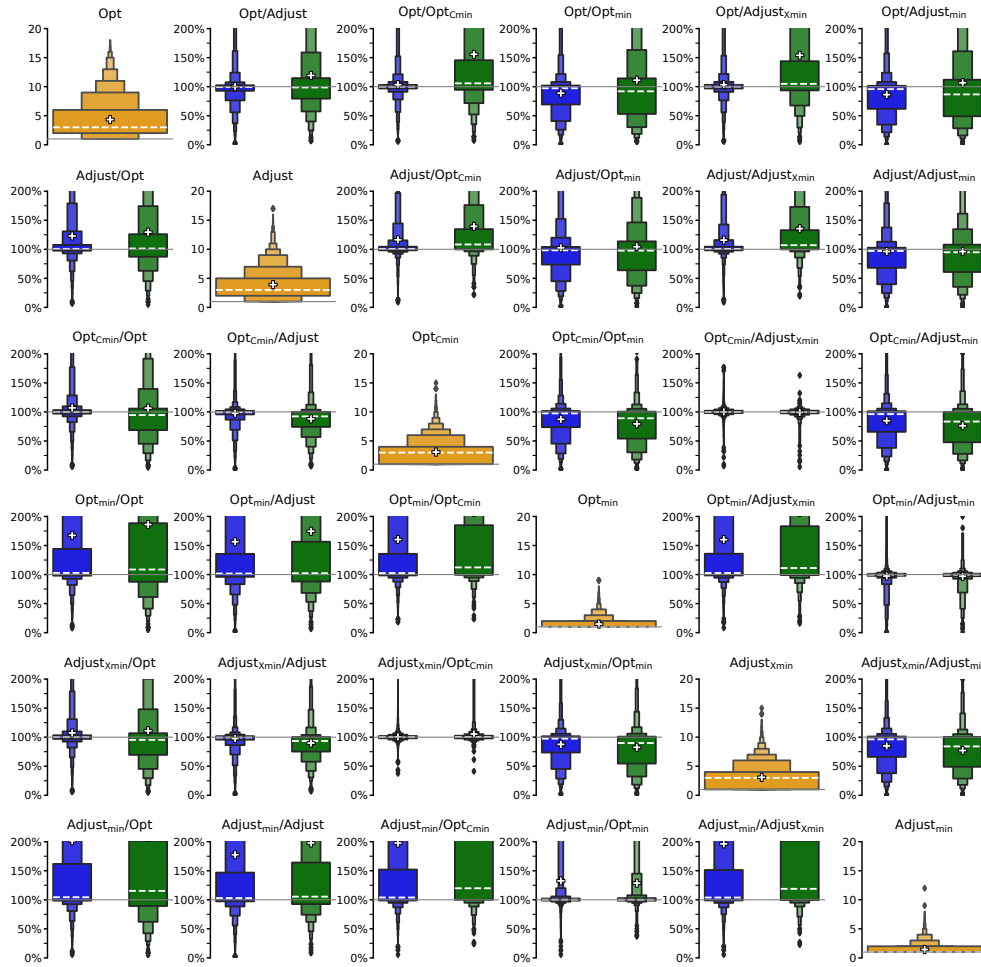


Figure S11: As in Fig. S2 but for MLP estimator, the nonlinear model, and $n = 1000$.

RF | random_lineargaussian | $n = 1000$ | $d_X = 1$
 93% of configurations are optimal | percentage of best-ranked among optimal / non-optimal configurations:
 Opt: 14 / 10 Adjust: 10 / 5 Opt_{Cmin}: 14 / 17 Opt_{min}: 24 / 24 Adjust_{xmin}: 13 / 19 Adjust_{min}: 22 / 22

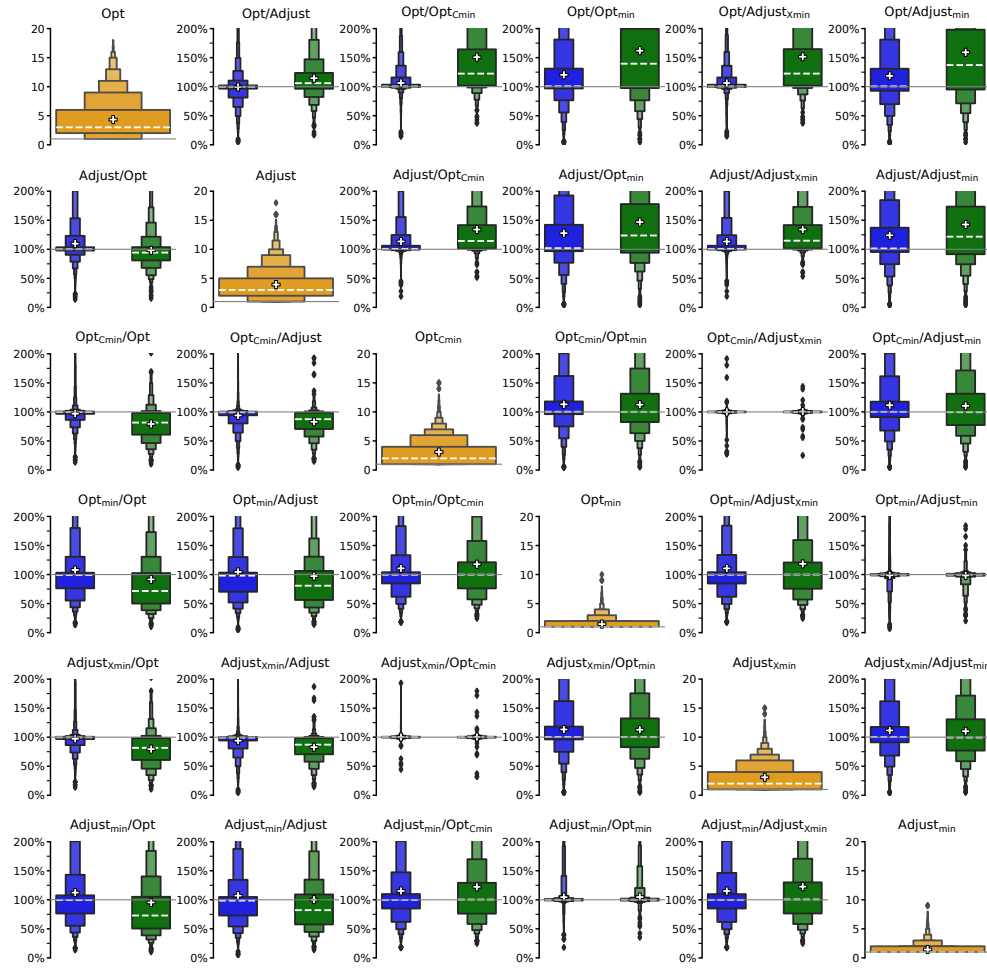


Figure S12: As in Fig. S2 but for RF estimator and $n = 1000$.

RF | random_nonlinearmixed | $n = 1000$ | $d_X = 1$
 93% of configurations are optimal | percentage of best-ranked among optimal / non-optimal configurations:
 Opt: 21 / 25 Adjust: 14 / 10 Opt_{Cmin}: 16 / 18 Opt_{min}: 15 / 12 Adjust_{xmin}: 16 / 23 Adjust_{min}: 14 / 10

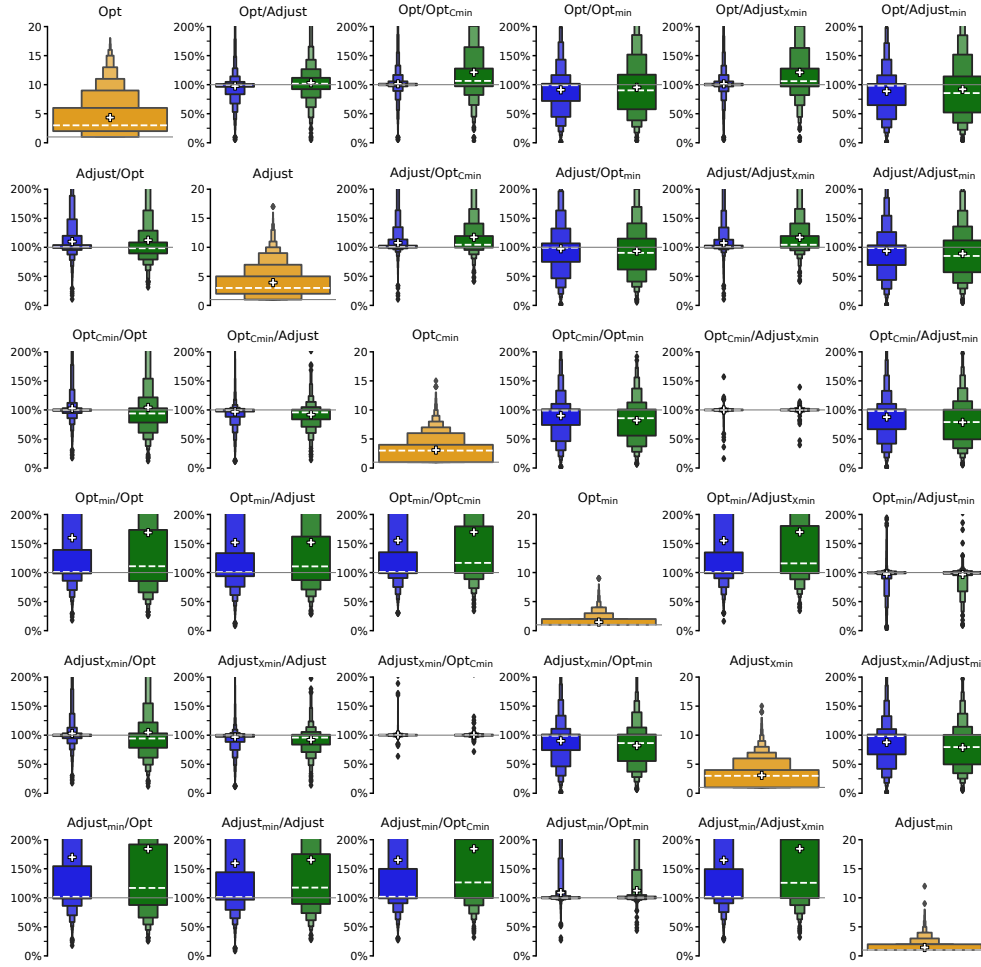


Figure S13: As in Fig. S2 but for RF estimator, the nonlinear model, and $n = 1000$.

DML | random_lineargaussian | $n = 1000$ | $d_X = 1$
 93% of configurations are optimal | percentage of best-ranked among optimal / non-optimal configurations:
 Opt: 24 / 19 Adjust: 13 / 4 Opt_{Cmin}: 20 / 29 Opt_{min}: 12 / 8 Adjust_{xmin}: 19 / 28 Adjust_{min}: 10 / 9

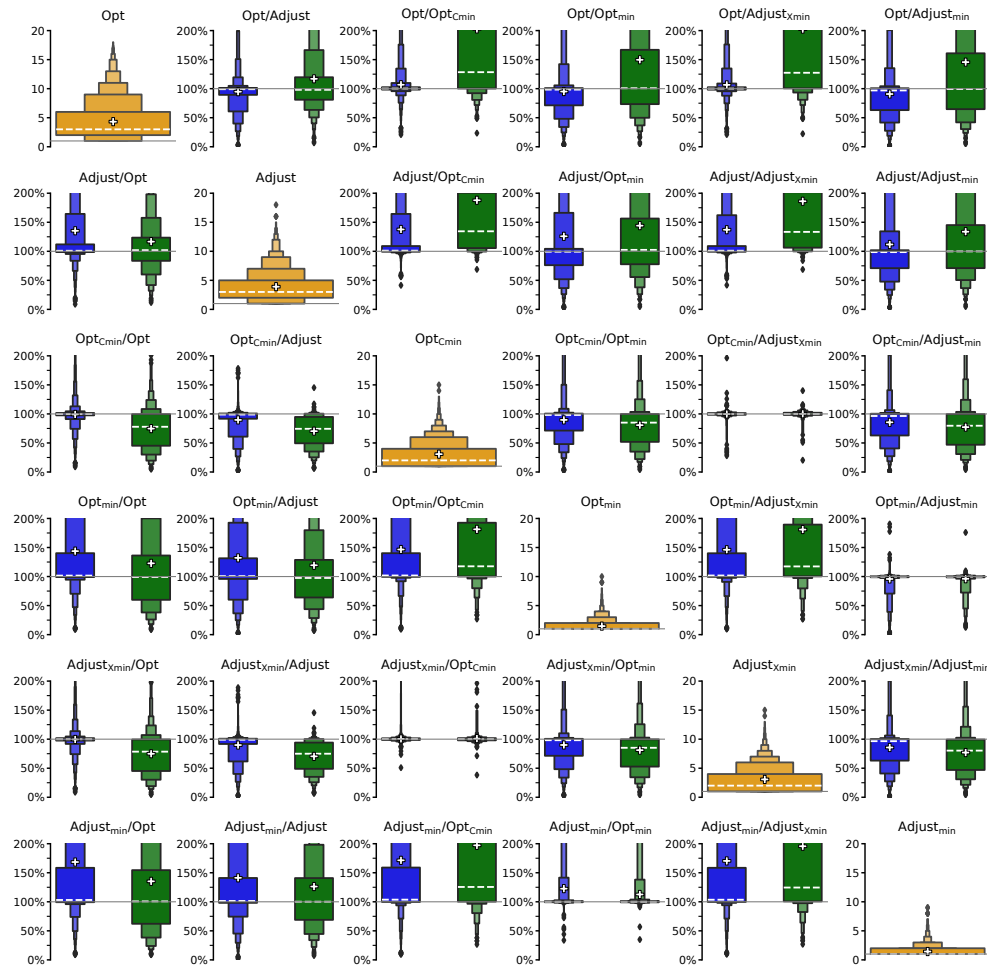


Figure S14: As in Fig. S2 but for DML estimator and $n = 1000$.

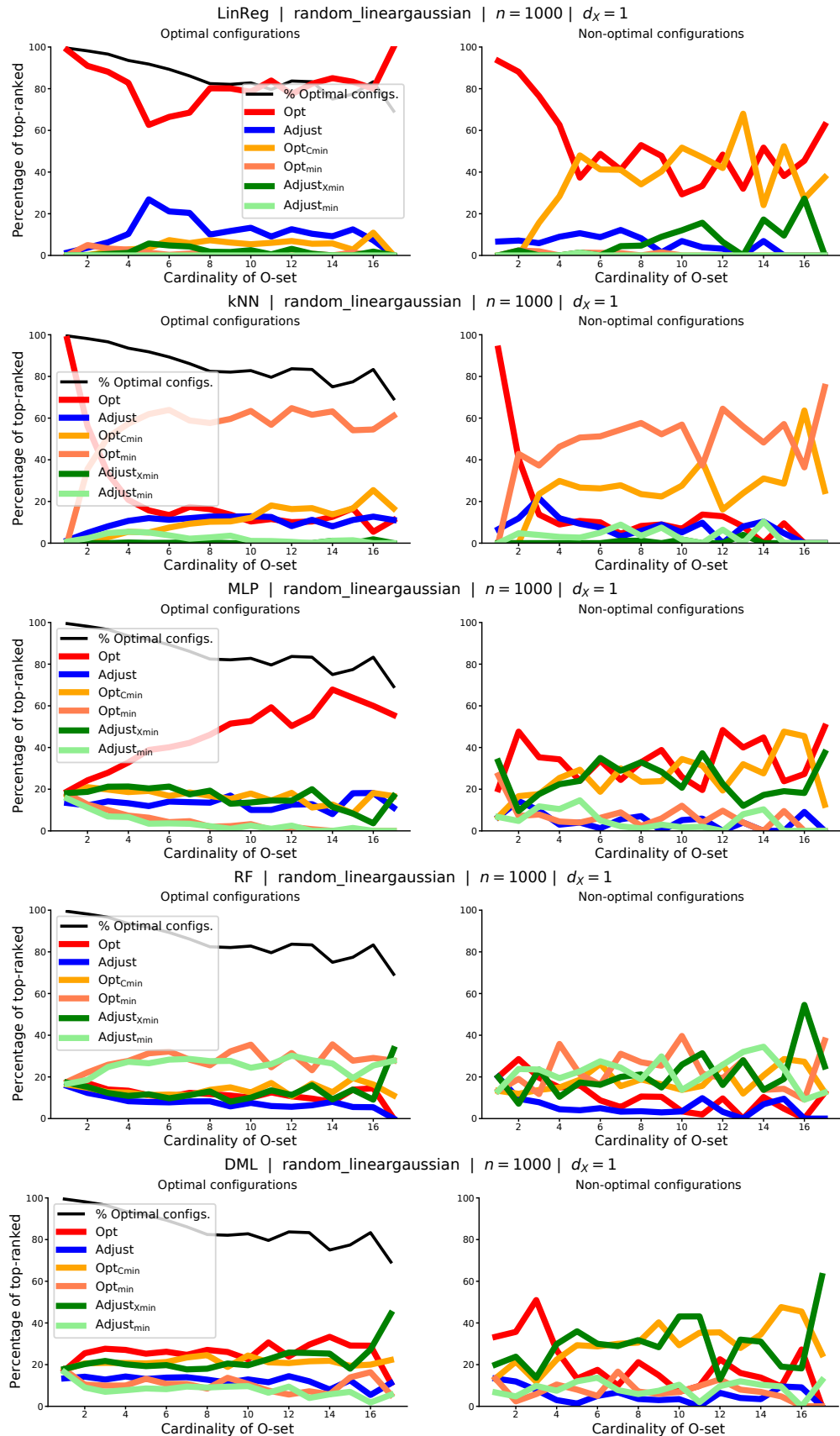


Figure S15: As in Fig. S7 but including non-parametric estimators for $n = 1000$.

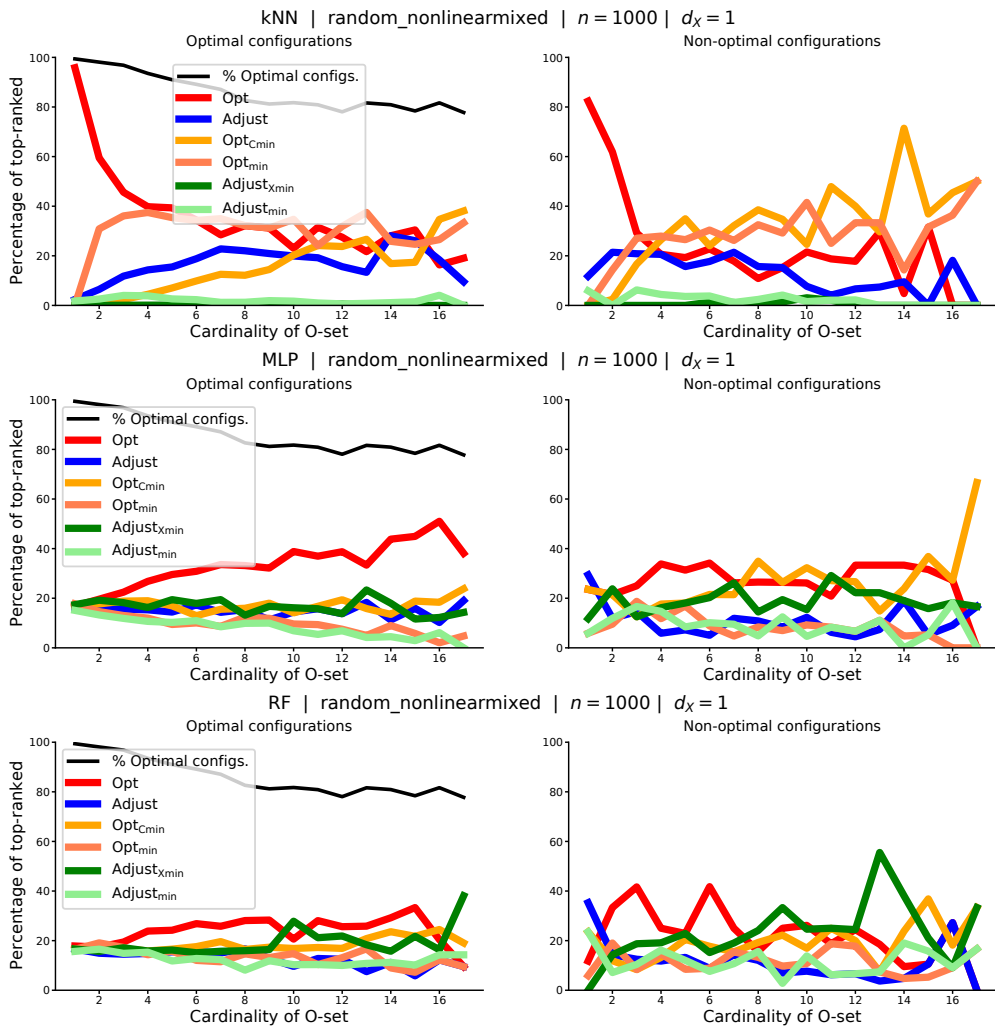


Figure S16: As in Fig. S15 but for nonlinear experiments.

References

- Philipp Bach, Victor Chernozhukov, Malte S. Kurz, and Martin Spindler. DoubleML – An object-oriented implementation of double machine learning in Python, 2021. arXiv:2104.03220 [stat.ML].
- Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21:C1–C68, 2018.
- Heike Hofmann, Hadley Wickham, and Karen Kafadar. Letter-value plots: Boxplots for large data. *Journal of Computational and Graphical Statistics*, 26(3):469–477, 2017.
- Judea Pearl. *Causality: Models, reasoning, and inference*. Cambridge University Press, 2009.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duchesnay. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12(Oct): 2825–2830, 2011. ISSN ISSN 1533-7928.
- Emilija Perković, Johannes Textor, and Markus Kalisch. Complete graphical characterization and construction of adjustment sets in markov equivalence classes of ancestral graphs. *Journal of Machine Learning Research*, 18:1–62, 2018.
- Thomas Richardson and Peter Spirtes. Ancestral graph markov models. *The Annals of Statistics*, 30(4):962–1030, 08 2002.
- Andrea Rotnitzky and Ezequiel Smucler. Efficient adjustment sets for population average treatment effect estimation in non-parametric causal graphical models. *arXiv preprint arXiv:1912.00306*, 2019.
- Benito van der Zander, Maciej Liśkiewicz, and Johannes Textor. Separators and adjustment sets in causal graphs: Complete criteria and an algorithmic framework. *Artificial Intelligence*, 270:1–40, 2019.
- Janine Witte, Leonard Henckel, Marloes H Maathuis, and Vanessa Didelez. On efficient adjustment in causal graphs. *Journal of Machine Learning Research*, 21(246):1–45, 2020.