



Swiss Federal Institute of Technology Zurich

Seminar for  
Statistics

Department of Mathematics

---

Semester Project

Fall 2019

---

Juan L Gamella

**Towards active ICP:  
experiment selection through stability**

---

Submission Date: February 16th 2020

---

Co-Adviser Dr. Christina Heinze-Deml  
Adviser: Prof. Dr. Nicolai Meinshausen

## Abstract

Causal DAG models cannot be fully identified using only observational data. Interventional data, that is, data originating from different experimental environments, improves identifiability; however, the improvement depends critically on the target and nature of the interventions carried out in each experiment. Since in real applications experiments tend to be costly, there is a need to perform the *right* interventions, in the sense of having to do as few of them as possible. Here we propose a new active learning (i.e. experiment selection) strategy based on causal invariance [Peters, Bühlmann, and Meinshausen 2016], which selects interventions that quickly reveal the direct causes of a target variable in the graph. We show that for general structural equation models, a direct cause appears on at least half of all stable sets [Pfister et al. 2019], and we further characterize the effect of interventions on such sets. We leverage these results to propose several intervention selection policies, which outperform a random policy in both population and finite-regime experiments.

## 1 Introduction

Causal models [Pearl et al. 2009] capture the causal relationships between variables and allow us, among other things, to make predictions of how a system behaves under interventions or distribution changes. In this sense, they are more powerful than probabilistic models, and can be seen as abstractions of more accurate mechanistic or physical models that retain enough power to answer interventional or counterfactual questions [Peters, Janzing, and Schölkopf 2017]. This allows them to maintain their predictive power in new, previously unseen environments [Pfister et al. 2019].

The question remains if for systems of interest such models can be learned directly from data. This problem is known in the literature as causal learning, and it is to causal models what statistical learning is to probabilistic models. Just as statistical learning, it suffers from the inherent difficulty of determining properties of a distribution from finite-sized samples. However, causal learning is additionally challenged by the fact that, even with full knowledge of the underlying distribution, some causal relationships cannot be established and causal models can generally not be fully identified from observational data alone [Pearl et al. 2009]. If discouraged, one should note that this limit applies to observational data, i.e. data that comes from a single environment or experimental setting. Furthermore, under additional assumptions about the model class and noise distributions, full identifiability is still possible [Peters et al. 2011]. In the general case, however, identifiability can only be improved by performing interventions (experiments). Examples of such interventions are abundant in the empirical sciences, from gene knockout experiments in biology to the assignment to treatment and control groups in clinical trials. For certain types of causal models there is an upper bound on the number of interventions required for full identifiability [Eberhardt 2008; Hauser and Bühlmann 2014].

A common way to express such relationships is through a directed acyclic graph (DAG), with variables as nodes and directed edges representing direct causal relationships [Spirtes et al. 2000]. This work deals with such models. In this case, the limit of identifiability means that, from observational data alone, the true graph cannot be distinguished from others that lie in the same Markov Equivalence class [Andersson, Madigan, Perlman, et al. 1997; Verma and Pearl 1990]. Therefore, if the full graph or properties of it (such as the direct causes of a target variable) are to be determined, there is a need to perform interventions. And, since experiments tend to be costly, there is a need to pick the *right* interventions, in the sense of having to do as few of them as possible.

### 1.1 Related work

Here we use the term *active causal learning* to refer to learning causal models from data while being able to actively perform interventions. In this setting, the goal is to sequentially improve identifiability, as opposed to the classical setting from machine learning [Settles 2009], where the goal is to sequentially increase prediction accuracy. Existing approaches can be said to fall broadly into two categories: Bayesian and non-Bayesian. The Bayesian approach, pioneered by the works of [Tong and Koller 2001; Murphy 2001], selects interventions which maximize a Bayesian utility function, generally the mutual information between the graph and the hypothetical sample that the experiment would produce. More recent works build on this approach by considering experiments performed in batches under budget constraints [Agrawal et al.

2019], or when expert knowledge is available [Masegosa and Moral 2013]. Non-Bayesian approaches employ graph-theoretic principles for selecting interventions that guarantee full identifiability [Eberhardt 2008, He and Geng 2008] and orient the maximum number of edges [Hauser and Bühlmann 2014; Ghassami et al. 2018].

Both approaches make different assumptions and suffer from different drawbacks. The Bayesian approach assumes a full probabilistic model, usually a joint Gaussian distribution. It is difficult to analyse how model misspecification influences the choice of experiments [Walker 2013]. Furthermore, it suffers from poor computational scaling [Ness et al. 2017] and several approximations have to be made even for small graphs [Agrawal et al. 2019]; this further complicates giving guarantees on the result. Graph theoretic approaches are agnostic to the underlying distribution, but they make two strong assumptions: (1) that the Markov equivalence class containing the true graph has been correctly identified, which is difficult with limited sample size, and (2) that interventions are perfectly informative.

## Invariant Causal Prediction

The work presented here is a first attempt at a new approach which falls into neither of the previous two categories. It is motivated by invariant causal prediction (ICP) [Peters, Bühlmann, and Meinshausen 2016], which allows recovering the direct causes of a variable (the *response*) from interventional data. The general idea is that the conditional distribution of a variable, given its direct causes, will remain invariant when intervening on arbitrary variables in the system other than itself. ICP considers the setting where different experimental conditions of a system exist and an independent, identically distributed sample of each environment is available. It then searches for sets of *plausible causal predictors*, namely sets of predictors which, if conditioned on, leave the distribution of the response invariant across the observed environments (the formal definition is given in section 3). Sets considered as plausible causal predictors given the current data are referred to as *accepted sets*, and the set of direct causes of the response (its parents in the graph) will be among them with high probability. ICP then returns the intersection of all accepted sets as an estimate of the direct causes.

While it does not retrieve the full graph, ICP presents some important advantages in the form of guarantees and more flexible assumptions. It requires neither knowledge of the Markov equivalence class, nor about the nature or location of the interventions performed in each environment, save that they must not be on the response. It assumes that the noise distribution of the response is independent from the direct causes and invariant across environments, but it makes no further assumptions on it or those of other variables; note that further assumptions can arise from the choice of tests for the invariance of the conditional distribution, but non-parametric tests can be chosen. Perhaps most importantly, ICP provides type-I-error control of the identified causes, namely that with high probability it will not mark as direct causes variables which are not. While this comes at a loss of power, this work shows that when coupled with an appropriate experiment selection strategy, ICP can quickly identify the direct causes while maintaining the aforementioned control.

## 1.2 Contributions

Using the framework of intervention stable sets [Pfister et al. 2019], we show that a direct cause appears on at least half of all stable sets (Proposition 2). We further characterize the effect of interventions on such sets, and use these results to construct intervention selection policies for ICP. Experiments are performed in both the population setting and the finite regime, comparing their performance to a random policy used as baseline. The proposed policies perform better in terms of the number of interventions required to correctly identify the causal predictors.

The framework used to generate synthetic interventional data, which allows generating and sampling from arbitrary structural equation models with arbitrary noise distributions, will be made available at [github.com/juangamella/sempler](https://github.com/juangamella/sempler).

## 1.3 Outline

Section 2 introduces the notions of prediction stability and stable sets of predictors [Pfister et al. 2019], and provides new theoretical results on how these sets are affected by interventions. Their relation to plausible causal predictors [Peters, Bühlmann, and Meinshausen 2016] is characterized in section 3. Section 4 leverages

these theoretical results to propose a first intervention selection policy, for which experimental results are provided in section 5. Section 6 contains the discussion and an outlook on future work.

## 2 How interventions affect the stable sets

We formalize in the following setting the assumptions required for the results derived in this section. The proofs can be found in the appendix.

**Setting 1** (adapted from setting 2 in [Pfister et al. 2019]) Let  $X \in \mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_p$  be predictor variables,  $Y \in \mathbb{R}$  a *response* variable and  $I = (I^1, \dots, I^m) \in \mathcal{I} = \mathcal{I}_1 \times \dots \times \mathcal{I}_m$  intervention variables which are unobserved and formalize the interventions present in the collection of intervention environments  $\mathcal{E}$ . Assume there exists a SEM  $\mathcal{S}^\mathcal{E}$  over  $(I, X, Y)$  that can be represented by a directed acyclic graph  $\mathcal{G}(\mathcal{S}^\mathcal{E})$ , in which the intervention variables are source nodes and do not appear in the right hand side of the assignment of  $Y$ , that is, assume there are no interventions on the response. For each  $e \in \mathcal{E}$  there is a SEM  $\mathcal{S}_e$  over  $(I_e, X_e, Y_e)$  such that  $\mathcal{G}(\mathcal{S}_e) = \mathcal{G}(\mathcal{S}^\mathcal{E})$ , in which only the equations with  $I_e$  on the right hand side change with respect to  $\mathcal{S}^\mathcal{E}$ . Furthermore, assume that the distribution of  $(I_e, X_e, Y_e)$  is absolutely continuous with respect to a product measure that factorizes.

Note that no assumption is made on the linearity of the structural equation models. In the active learning setting, at each iteration an intervention target is selected and a sample is collected from the new experimental environment. No further assumptions are made on the size or type of the intervention. Denote by  $\mathcal{E}_t = \{e_i : i \in \{1, \dots, t\}\}$  the set of observed environments up to iteration  $t$  of the algorithm, and assume  $\mathcal{E}_t \subseteq \mathcal{E}_{t+1}$ . As for simplifying notation, let  $\text{DE}(S) = \{j \in \{1, \dots, p\} \mid \exists i \in S : j \in \text{DE}(i)\}$  denote the descendants of variables in a set  $S$ , let  $\text{PA}(i)$  be the parents of  $X_i$  and let  $\text{PA}(S) = \{j \in \{1, \dots, p\} \mid \exists i \in S : j \in \text{PA}(i)\}$  denote the parents of variables in a set  $S$ . Note that the descendants of a variable include the variable itself, i.e.  $i \in \text{DE}(i)$ .

### 2.1 Stable sets and the stable blanket

The notion of intervention stable sets is introduced in [Pfister et al. 2019], and allows characterizing sets of plausible causal predictors from d-separation statements in the graph. While stable sets are not generally equivalent to the sets of plausible causal predictors accepted by ICP, the approach here is to derive theoretical results for them, and then see under which conditions these apply to the accepted sets (section 3).

**Definition 2.1** (intervention stable set Pfister et al. 2019). *Given setting 1 and a set of environments  $\mathcal{E}$ , we call a set  $S \subseteq \{1, \dots, p\}$  intervention stable under  $\mathcal{E}$  if the d-separation  $I \perp\!\!\!\perp_{\mathcal{G}} Y \mid S$  holds in  $\mathcal{G}(S^*)$  for any intervention  $I$  which is active in an environment  $e \in \mathcal{E}$ .*

In other words, a set of predictors is stable if it d-separates the response from all interventions. From now on, let  $\mathbb{S}_{\mathcal{E}}$  denote the collection of sets which are intervention stable under  $\mathcal{E}$ .

The Markov blanket is the smallest set of variables which d-separate the response from all other variables [Pearl 1988], and is therefore optimal in terms of prediction accuracy. The stable blanket is the subset of the Markov blanket which is intervention stable given the present interventions.

**Definition 2.2** (stable blanket Pfister et al. 2019). *Define the set*

$$N_t^{\text{int}} := \{1, \dots, p\} \setminus \{j \in \{1, \dots, p\} \mid \exists k \in CH_t^{\text{int}}(Y) : j \in \text{DE}(X_k)\},$$

where  $CH_t^{\text{int}}(Y)$  are the children of  $Y$  which have been directly intervened in at least one environment  $e \in \mathcal{E}_t$ , and  $\text{DE}(X_k)$  are the descendants of  $X_k$ , itself included. Then, the stable blanket, denoted as  $SB_t(Y)$  is defined as the smallest set  $S \subseteq N_t^{\text{int}}$  such that

$$X_j \perp\!\!\!\perp_{\mathcal{G}} Y \mid S \quad \forall j \in N_t^{\text{int}} \setminus S,$$

where  $\perp\!\!\!\perp_{\mathcal{G}}$  denotes d-separation in the graph  $\mathcal{G}(\mathcal{S}^{\mathcal{E}_t})$ .

Note that  $N_t^{\text{int}}$  is the set of all variables, excluding the directly intervened children of  $Y$  and their descendants. Furthermore, note that our assumption  $\mathcal{E}_t \subseteq \mathcal{E}_{t+1}$  implies that  $\text{CH}_t^{\text{int}}(Y) \subseteq \text{CH}_{t+1}^{\text{int}}(Y)$ . A more intuitive characterization of the stable blanket follows.

**Theorem 1** (stable blanket Pfister et al. 2019). *The stable blanket consists of the parents of  $Y$ , the children of  $Y$  in  $N_t^{\text{int}}$  and the parents of such children.*

From now on, let  $\text{SB}_t(Y)$  denote the the stable blanket of the response under environment  $\mathcal{E}_t$ , i.e. the stable blanket at iteration  $t$  of the active learning procedure. If there are no interventions, the stable blanket is equal to the Markov blanket; if there are sufficiently many informative interventions, the stable blanket contains only the parents, or direct causes, of the response. The following proposition formalizes what an informative intervention is in this sense, and is the first main result of this section: The stable blanket is only affected by direct interventions on children which have not yet been directly intervened on. Furthermore, the change in the stable blanket consists in the removal of

1. the descendants of such children, and
2. variables which are not children of the response and whose children are all among such descendants.

**Proposition 1** (effect of interventions on the stable blanket). *Let  $\mathcal{E}_t, \mathcal{E}_{t+1}$  be sets of observed environments such that  $\mathcal{E}_t \subseteq \mathcal{E}_{t+1}$ . Let  $I_{t+1}$  be the variables intervened in  $e_{t+1}$ , and let  $S = \text{CH}_{t+1}^{\text{int}}(Y) \setminus \text{CH}_t^{\text{int}}(Y) = I_{t+1} \cap \text{CH}(Y) \cap N_t^{\text{int}}$  be the children of the response intervened on for the first time in  $e_{t+1}$ . Then,*

$$\begin{aligned} \text{SB}_t(Y) = & \text{SB}_{t+1}(Y) \cup \\ & (DE(S) \cap \text{SB}_t(Y)) \cup \\ & \{j \in \text{SB}_t(Y) : j \notin DE(S) \wedge j \notin (\text{CH}(Y) \cap N_t^{\text{int}}) \wedge \text{CH}(j) \subseteq DE(S)\}. \end{aligned}$$

In consequence, interventions on other variables in the stable blanket, such as parents of the response or parents of children of the response, do not modify the stable blanket. They do, however, modify the set of intervention stable sets, and this can be used to infer certain information about the graph structure and guide an active learning policy.

## 2.2 How do interventions affect the intervention stable sets?

If a variable upstream of the response is intervened, at least one variable in the directed path from it to the response must be present in all stable sets (Lemma 1). If the intervened variable is a parent, it appears on all stable sets (Lemma 2).

**Lemma 1** (interventions on upstream variables). *Let  $\mathcal{E}$  be a set of observed environments and let  $j \in \text{AN}(Y)$  such that it is directly intervened in  $\mathcal{E}$ . Then,*

$$S \subseteq \{1, \dots, p\} \text{ is intervention stable} \implies S \cap DE(j) \cap \text{AN}(Y) \neq \emptyset.$$

**Lemma 2** (intervened parents appear on all intervention stable sets). *Let  $\mathcal{E}$  be a set of observed environments and let  $j \in \text{PA}(Y)$  such that it is directly intervened in  $\mathcal{E}$ . Then,*

$$S \subseteq \{1, \dots, p\} \text{ is intervention stable} \implies j \in S.$$

The opposite is true when intervening on children of the response, namely,

**Lemma 3** (sets containing descendants of directly intervened children are unstable). *Let  $i \in \text{CH}^{\text{int}}(Y)$ . Then, any set  $S \subseteq \{1, \dots, d\}$  which contains descendants of  $i$  is not intervention stable.*

Changes in the stable blanket yield additional information about the graph structure. Proposition 1 tells us that, after an intervention, variables that drop out of the stable blanket are either children or parents of children. The following lemma allows us to classify the variables that drop out of the stable blanket.

**Lemma 4** (stability of parents of intervened children). *Let  $\mathcal{E}_t, \mathcal{E}_{t+1}$  be sets of observed environments such that  $\mathcal{E}_t \subseteq \mathcal{E}_{t+1}$ . Assume an intervention occurred such that  $j \in SB_t(Y)$  but  $j \notin SB_{t+1}(Y)$ , and let  $S = CH_{t+1}^{int}(Y) \setminus CH_t^{int}(Y)$  be the set of newly intervened children. Then, if  $SB_{t+1}(Y) \cup \{j\}$  is intervention stable,*

$$j \in \{j \in SB_t(Y) : j \notin DE(S) \wedge j \notin (CH(Y) \cap N_t^{int}) \wedge CH(j) \subseteq DE(S)\}.$$

In other words, if a variable drops out of the stable blanket, but is stable together with the new stable blanket, then (1) it is not a descendant of the response or intervened children, and (2) all its children in the blanket have been directly intervened on; this allows us to identify parents of children of the response which satisfy (1) and (2). Additional structure can be inferred by considering the number of stable sets in which a predictor appears:

**Definition 2.3** (stability ratio). *Given a set of environments  $\mathcal{E}$ , the stability ratio of a variable  $i \in \{1, \dots, p\}$  is defined as*

$$r_{\mathcal{E}}(i) := \frac{1}{|\mathbb{S}_{\mathcal{E}}|} \sum_{S \in \mathbb{S}_{\mathcal{E}}} \mathbb{1}\{i \in S\},$$

*i.e. in which proportion it appears in the intervention stable sets under  $\mathcal{E}$ .*

**Corollary 1.1.** *From Lemma 2 and Lemma 3 it follows that*

1. *parents which are directly intervened in at least one environment in  $\mathcal{E}$  have a ratio of 1, and*
2. *descendants of children directly intervened in at least one environment in  $\mathcal{E}$  have a ratio of 0.*

The following proposition is the other main result of this section, and is directly used in section 4 to construct an active learning policy.

**Proposition 2** (parents appear on at least half of all stable sets). *Let  $\mathcal{E}$  be any set of observed environments. Then, for any  $j \in \{1, \dots, p\}$ ,*

$$r_{\mathcal{E}}(j) < 1/2 \implies j \notin PA(Y).$$

### 3 From stable sets to causal predictors

The results derived in section 2 apply to intervention stable sets; if we are to use these results to construct an active learning policy for ICP, we need to know under which conditions they apply directly to the sets of plausible causal predictors.

**Definition 3.1** (plausible causal predictors Peters, Bühlmann, and Meinshausen 2016). *We call a set of variables  $S \subseteq \{1, \dots, p\}$  plausible causal predictors under a set of environments  $\mathcal{E}$  if there exists a vector of coefficients  $\beta \in \mathbb{R}^p$  such that  $\beta_i = 0$  if  $i \notin S$  and*

$$\text{for all } e \in \mathcal{E}, X^e \text{ has arbitrary distribution and} \tag{1}$$

$$Y^e = \mu + X^e \beta + \varepsilon^e, \quad \varepsilon^e \sim F_{\mathcal{E}} \text{ and } \varepsilon^e \perp\!\!\!\perp X_S^e, \tag{2}$$

*where  $\mu \in \mathbb{R}$  is an intercept term,  $\varepsilon^e$  is random noise with zero mean, finite variance and the same distribution  $F_{\mathcal{E}}$  across all  $e \in \mathcal{E}$ . Let  $\mathbb{C}_{\mathcal{E}}$  denote the collection of sets which are plausible causal predictors under  $\mathcal{E}$ .*

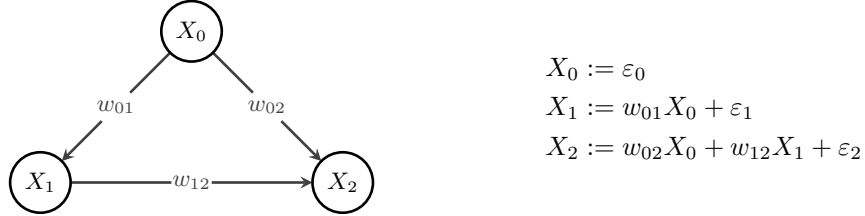
**Corollary 1.2.** *Let  $\mathcal{E}_t, \mathcal{E}_{t+1}$  be sets of observed environments such that  $\mathcal{E}_t \subseteq \mathcal{E}_{t+1}$ . Then, it follows that if  $S$  is not a set of plausible causal predictors under  $\mathcal{E}_t$ , it also is not under  $\mathcal{E}_{t+1}$ .*

In an active learning setting, where a new environment is added in each iteration, corollary 1.2 provides a speed up, as only the sets accepted in the previous iteration have to be considered. Given a collection of environments  $\mathcal{E}$ , the collection of accepted sets of the ICP algorithm (Peters, Bühlmann, and Meinshausen 2016) is an estimate of  $\mathbb{C}_{\mathcal{E}}$ . The following proposition establishes the relationship between intervention stable sets and sets of plausible predictors.

**Proposition 3** (intervention stable sets are plausible causal predictors). *Let  $\mathcal{E}$  be a set of observed environments. Then, for all intervention stable sets  $S \subseteq \{1, \dots, p\}$ , it holds that  $S \in \mathbb{C}_{\mathcal{E}}$ .*

While  $\mathbb{S}_{\mathcal{E}} \subseteq \mathbb{C}_{\mathcal{E}}$ , it is important to note that, even under the faithfulness assumption, it is not generally true that  $\mathbb{S}_{\mathcal{E}} = \mathbb{C}_{\mathcal{E}}$ . An example follows.

**Example 3.1.** Take the following SEM,



with  $\varepsilon_i$  noise variables such that  $\mathbb{E}[\varepsilon_i] =: \mu_i$ ,  $\text{Var}(\varepsilon_i) =: \sigma_i^2$  and  $\varepsilon_i \perp\!\!\!\perp \varepsilon_j \forall i, j$ . Performing an OLS regression of  $Y := X_1$  on the predictor sets  $S_0 = \{0\}$  and  $S_2 = \{2\}$  yields the following coefficients and intercepts:

$$\begin{aligned} \beta^{S_0} &= w_{01}, \quad \mu^{S_0} = \mu_1, \\ \beta^{S_2} &= \frac{\sigma_0^2(w_{01}^2w_{12} + w_{01}w_{02}) + \sigma_1^2w_{12}}{\sigma_0^2(w_{01}w_{12} + w_{01}w_{02})^2 + \sigma_1^2w_{12} + \sigma_2^2}, \quad \mu^{S_2} = \mu_1 + w_{01}\mu_0 - \beta^{S_2}(\mu_2 + w_{12}\mu_1 + \mu_0(w_{01}w_{12} + w_{02})). \end{aligned}$$

If we additionally assume  $\mu_i = 0$ ,  $w_{ij} = 1 \forall i, j$  and  $\sigma_1^2 = \sigma_2^2 = 1$ , the above expressions become

$$\begin{aligned} \beta^{S_0} &= 1, \quad \mu^{S_0} = 0, \\ \beta^{S_2} &= \frac{2\sigma_0^2 + 1}{4\sigma_0^2 + 2} = \frac{1}{2}, \quad \mu^{S_2} = 0. \end{aligned}$$

Consider now an intervention on  $X_0$ . We have that  $S_0 = \{0\}$  is intervention stable and a set of plausible causal predictors.  $S_2 = \{2\}$  does not d-separate  $Y$  from the intervention on  $X_0$ , and is not intervention stable; however, for interventions that affect only the variance of  $X_0$  (i.e.  $\sigma_0^2$ ),  $S_2$  is a set of plausible causal predictors. Under this setting, we have that  $\mathbb{S}_{\mathcal{E}} \subset \mathbb{C}_{\mathcal{E}}$ .

In the above, note that the regression on  $X_2$  yields a lower mean squared error. Furthermore, in the population setting the OLS regression over  $X_1$  and  $X_2$  sets the coefficient for  $X_1$  to zero. This highlights several facts (in the population setting):

1. The stable blanket can be defined as the smallest intervention stable set with minimum MSE, but without further assumptions this does not hold sets of plausible causal predictors.
2. While OLS sets the coefficients of variables outside the Markov blanket to zero, variables inside the Markov blanket might also have zero coefficient. This means that taking the non-zero coefficients of OLS as an estimate of the Markov blanket can yield subsets of it.

While in theory this means that  $\mathbb{S}_{\mathcal{E}} \neq \mathbb{C}_{\mathcal{E}}$ , one might ask how often this happens in practice. In the above example, this only happens when we set the weights, means and variances to very particular values. When these parameters are sampled from a continuous distribution, one might argue that the set of parameters for which  $\mathbb{S}_{\mathcal{E}} \neq \mathbb{C}_{\mathcal{E}}$  has probability zero. This is formalized in the following conjecture,

**Conjecture 3.1.** *Assume setting 1 with a linear SEM, and let  $\theta \in \Theta$  be its parameters, i.e. weights, intercepts and noise variances. Let  $\Pi$  be an absolutely continuous distribution over  $\Theta$ . Then  $\mathbb{S}_{\mathcal{E}} = \mathbb{C}_{\mathcal{E}}$   $\Pi$ -a.s.*

## 4 Constructing an active learning policy

Even in the population setting, the capacity of ICP to retrieve the parents relies heavily on the informativeness of the environments. For example, if none of the interventions are upstream of the response, the empty set is intervention stable and is returned as estimate of the parents.

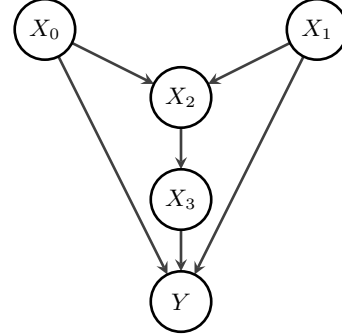
While [Peters, Bühlmann, and Meinshausen 2016] gives some sufficient conditions for the identifiability of the true causal predictors, it is not entirely clear what an optimal intervention is. If we make the assumption that  $\mathbb{C}_{\mathcal{E}} = \mathbb{S}_{\mathcal{E}}$ , by Lemma 2 we know that a direct intervention on a parent is sufficient for it to appear in the ICP estimate. However, it is not a necessary condition, as shown in the following example.

**Example 4.1.** Let  $\mathcal{E}$  be a collection of two environments: one without interventions and one with a direct intervention on  $X_2$ . The intervention stable sets are

$$\begin{aligned}\mathbb{S}_{\mathcal{E}} = & \{0, 1, 2\}, \\ & \{0, 1, 3\}, \\ & \{0, 1, 2, 3\}.\end{aligned}$$

Therefore,

$$S(\mathcal{E}) = \bigcap_{S: S \in \mathbb{S}_{\mathcal{E}}} S = \{0, 1\},$$



which shows that parents can appear in the intersection of intervention stable sets without being directly intervened on.

In the above example, a direct intervention on  $X_2$  is very informative, as it reveals two parents simultaneously. To the best of our knowledge it is not clear when situations like the above arise, or how they can be detected from the accepted sets. Therefore, direct interventions on the parents are treated as “maximally informative”, and the goal of the proposed policies is to pick such interventions.

The implemented policies select the next intervention target by looking at the sets accepted by ICP in the previous iteration. Since these are not available when selecting the initial intervention, the sample from the initial environment can be used to guide a first choice. For now, only single-variable, shift interventions are considered. By corollary 1.2, at each iteration ICP can consider only the sets accepted in the previous iteration, providing a substantial speed up. An outline of the resulting active learning procedure is provided in Algorithm 1.

---

### Algorithm 1: Active ICP

---

**Output:**  $S(\mathcal{E})$  estimate of the parents of the response

**Input :** `policy` an intervention selection policy,  
 $(X^0, Y^0)$  sample from initial environment,  
 $T$  number of iterations

$\mathcal{E}_0 \leftarrow \{(X^0, Y^0)\};$

`accepted_sets`  $\leftarrow$  all sets of predictors;

`next_intervention`  $\leftarrow$  `policy.first_intervention`( $\mathcal{E}_0$ );

**for**  $t = 1 : T$  **do**

    perform `next_intervention` and collect sample  $(X^t, Y^t)$ ;

$\mathcal{E}_t \leftarrow \mathcal{E}_{t-1} \cup \{(X^t, Y^t)\};$

`accepted_sets`,  $S(\mathcal{E}_t) \leftarrow$  ICP( $\mathcal{E}_t$ , `accepted_sets`) ; // see corollary 1.2

`next_intervention`  $\leftarrow$  `policy.next_intervention`(`accepted_sets`);

**end**

**return**  $S(\mathcal{E}_T)$

---



## Proposed policies

To increase the chances of picking a parent of the response as an intervention target, the proposed policies can make use of three heuristics:

1. (*markov heuristic*) Selecting intervention targets from within the Markov blanket, which contains the parents. Under linearity and some assumptions (see example 3.1), in the population setting the Markov blanket can be directly obtained from an OLS regression over all predictors. In the finite regime, we turn to the Lasso [Tibshirani 1996] to obtain an estimate.
2. (*empty-set heuristic*) If after an intervention the empty set is accepted, the distribution of the response did not change as a result of the intervention. This means that the target is not a direct cause of the response or, under faithfulness, that it is not upstream of the response. In any case, we discard the target from future interventions.
3. (*ratio heuristic*) By Proposition 2, if a variable appears on less than half of all accepted sets, it is not a parent; therefore, we do not add it to the pool of possible intervention targets for the current iteration. Note that unlike in (2.), we do not discard it from future interventions. This is important in the finite regime, where parents may for some iterations appear in less than half of all accepted sets.

In the population setting, making use of all three heuristics reduces the number of interventions required to identify all direct causes. In the finite regime more care has to be taken; in the experiments we evaluate the effect of each heuristic individually, and show that relying on an estimate of the Markov blanket (1.) has a detrimental effect on the performance of the policy.

## 5 Experiments

We evaluate policies that use different combinations of the heuristics in both the population and finite sample setting, using simulated data from randomly chosen linear SEMs.

### 5.1 Population setting

Experiments are first run the population setting, as this simplifies an initial evaluation of the proposed policy, while illustrating the challenges faced in the finite regime:

1. Since interventions are perfectly informative, the performance of the policies can be compared exclusively in terms of their choice of targets, without worrying about (1) the parameters of the intervention, and (2) how many samples must be allocated to the experiment, neither of which are trivial problems.
2. We can ignore the problem of estimating the Markov blanket, which is simply taken as the variables with non-zero coefficient in an OLS regression over all predictors. See example 3.1 for a situation where this approach causes problems.
3. By Lemma 2 we have that in the population setting at most  $p$  different interventions are needed to produce the correct estimate,  $p$  being the number of predictors. This directly yields a limit on the number of iterations for which the active procedure has to be run.

In this setting, the proposed policy makes use of all three heuristics. We evaluate its performance against two baseline policies of increasing complexity. The first picks intervention targets at random from the predictor variables. The second selects variables to intervene on at random from the Markov blanket of the response, as this increases the chances of picking a parent. Note that in the cases where the Markov blanket estimate contains the whole graph, the two baseline policies are equivalent. Furthermore, when the Markov blanket estimate is composed exclusively by the parents, the second baseline and the proposed policy are equivalent.

For the experiments, 1000 linear structural equation models of size 15 are randomly generated, and the response is selected so the Markov blanket contains more than just the parents. The weights, intercepts and noise variances are sampled uniformly at random from  $[0, 1]$ . In the population setting no further assumptions

are made on the noise distributions, besides having finite mean and variance to perform the OLS regression. Our experiments with SEMs of different size and parameters yielded very similar results, and are not shown separately. For every SEM, each policy is run 32 times with different random seeds, as the policies have a stochastic component (Figure 3); they differ in how narrowly they constrain the possible intervention targets, but they choose randomly among them. We compare how quickly (in terms of interventions) the policies recover the causal predictors (Figure 1), and how many interventions they require in total (Figure 2). In both measures, the proposed policy performs better than both baselines.

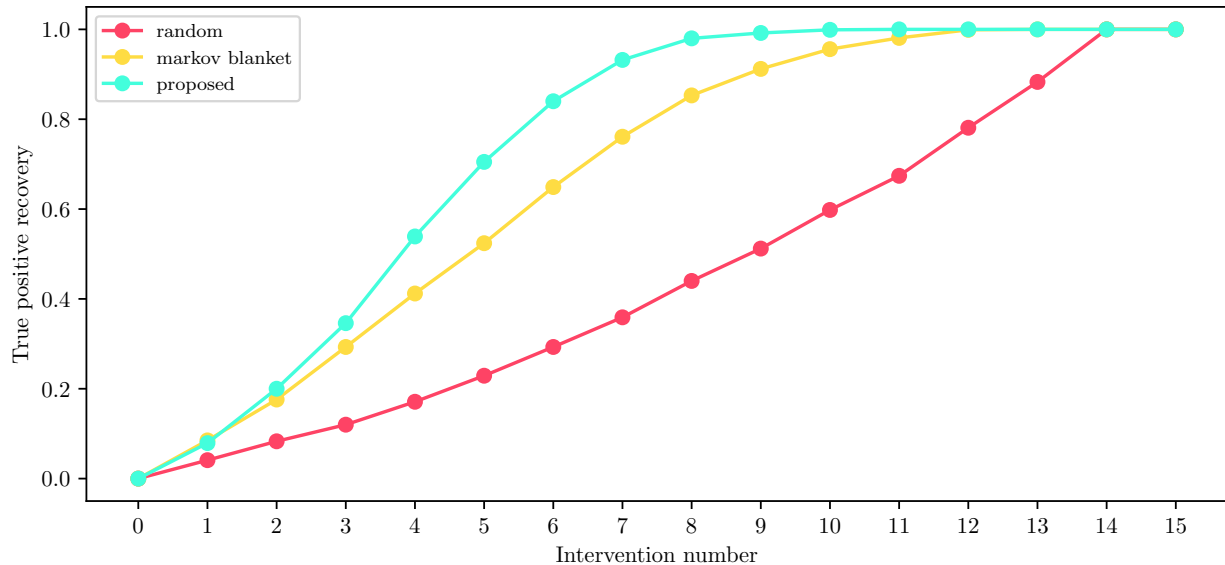


Figure 1: (population setting) Proportion of SEMs for which all direct causes have been recovered as a function of the number of interventions. Note that in the first intervention the proposed policy and the Markov blanket policy perform equally, as they decide among the same pool of intervention targets (the markov blanket).

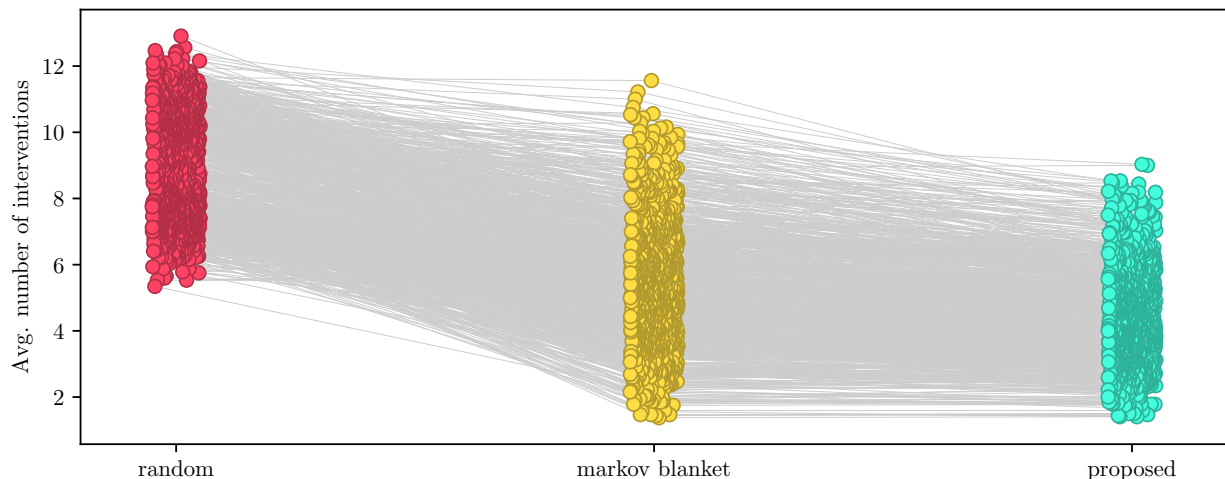


Figure 2: (population setting) Average number of interventions until the causal predictors are identified, for each one of the 1000 SEMs. Each SEM is represented by a dot and connected across policies by a grey line. On average, the random policy employs 8.9 interventions to produce the correct estimate, while the Markov blanket policy requires 5.5 and the proposed policy 4.3.

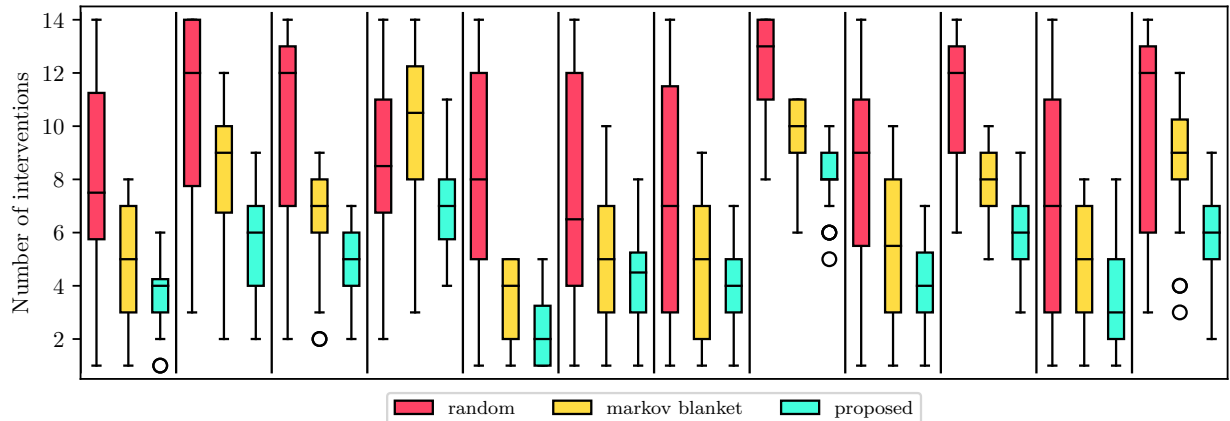


Figure 3: (population setting) Box plots of the number of interventions employed by each policy, for 12 SEMs selected at random from the total 1000. In this case, the policies are run a total of 100 times to better characterize their behaviour. By narrowing the choice of possible intervention targets, the Markov blanket policy and the proposed policy reduce both the average and maximum number of interventions they employ to identify the causal predictors. The difference in the minimum number of interventions between the shown SEMs is due to the number of parents of the response.

## 5.2 Finite sample setting

For the experiments in the finite regime, 300 linear structural equation models of size 8 are randomly generated. The remaining parameters are sampled as in the population setting. To simplify the implementation, we assume that the underlying noise distributions are Gaussian, and set ICP to use a two-sample t-test and F-test to check the invariance of the conditional distribution of the response. It is important to note that this is not a necessary requirement: the results derived in section 2 (e.g. Proposition 2) apply to arbitrary SEMs with arbitrary noise distributions, and ICP can use other statistical tests, including non-parametric ones. However, we expect that the effect of the number of samples on the results will be different under different noise distributions and tests.

In this setting, we individually evaluate the effect of the three heuristics put forward in section 4, as well as combinations of them. In total we have 7 policies, each using a different combination of heuristics, plus the random policy used as baseline. For the sample allocation, we fix the size of the sample collected per intervention; interventions are shift interventions with mean 10 and variance 1. We perform experiments for 10, 100 and 1000 observations per sample. Like in the population setting, the policies are run a total of 32 times with different random seeds. The maximum number of iterations is set to 50. We again compare the policies in terms of how quickly they recover the causal predictors (Figure 4), and the total number of interventions employed (Figure 5). Even at the smallest sample size, the policies that employ the *empty-set* and *ratio* heuristics outperform all others, including the random baseline.

The results show some interesting patterns. Relying on the Markov blanket estimate (policies labelled with “markov”) leads to a poor overall performance, independently of what other heuristics are used. In the experiments this estimate is obtained by performing an L1-regularized least squares regression on all predictors (i.e. the Lasso [Tibshirani 1996]); the lambda parameter is picked for each SEM by cross validation. This estimate generally contains some of the parents, which allows the policies to quickly identify them in the first few iterations, as can be seen in Figure 4; however, after these parents are identified, if not all parents are contained in the estimate the policies become stuck performing non-informative interventions. As a result, for many of the generated SEMs, the correct estimate is not produced after reaching the maximum number of iterations. As expected, this effect is minimized at larger sample sizes, where the Lasso yields a better estimate of the Markov blanket.

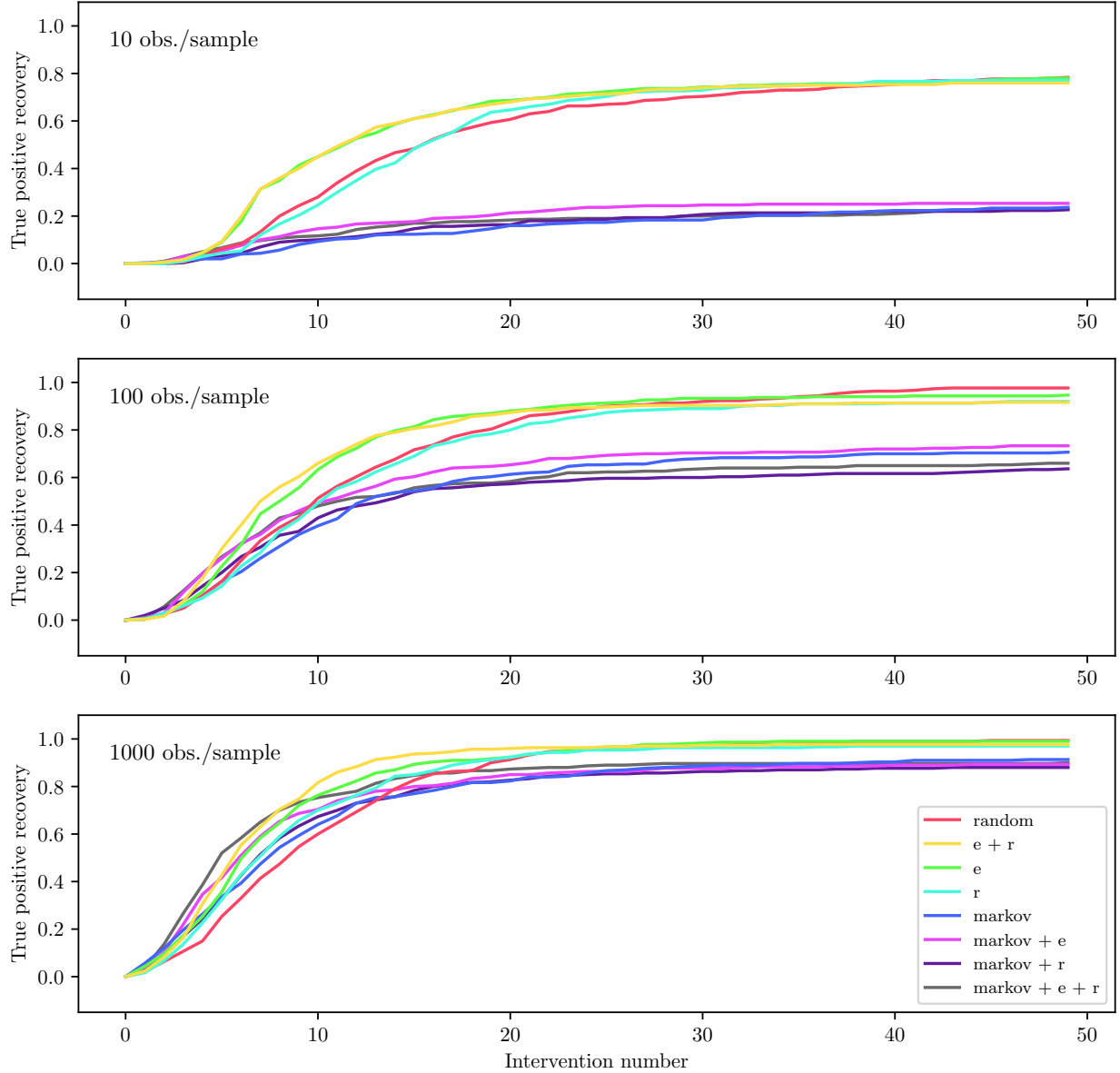


Figure 4: (finite regime) Proportion of SEMs for which all direct causes have been recovered as a function of the number of interventions. Policies are labelled according to the heuristics they employ: selecting intervention targets only from the Markov blanket estimate (markov); discarding targets for which, after being intervened, the empty set was accepted (e), and picking targets which appear on at least half of the accepted sets (r). Note that the performance of the *ratio* heuristic depends strongly on the number of samples. This is what causes the boost in performance of the yellow policy at 1000 obs./sample, which combines it with the *empty-set* heuristic, which performs well at all sample sizes.

The other two heuristics, i.e. the *ratio* and *empty-set*, yield an improvement over the random policy, and together produce the best-performing policy. That said, a point must be made about consistency in the number of interventions. For a fixed sample size, only the random policy and one which relies exclusively on the *ratio* heuristic are consistent in the number of interventions. Note that the *ratio* heuristic discards targets only for the current iteration, not future ones; this way, parents which momentarily appear on less than half of all accepted sets are not prevented from being intervened on in the future. From the previous paragraph it is clear why policies that rely on the Markov blanket estimate are not consistent in the number

of interventions. In the case of the *empty-set* heuristic, with some probability the empty set will be wrongly accepted after an intervention on a parent, which is then discarded from future interventions. This is not solved by lowering the level of the tests. For large sample sizes the problem is attenuated by the fact that the intervention can already be enough for ICP to identify the target as a direct cause.

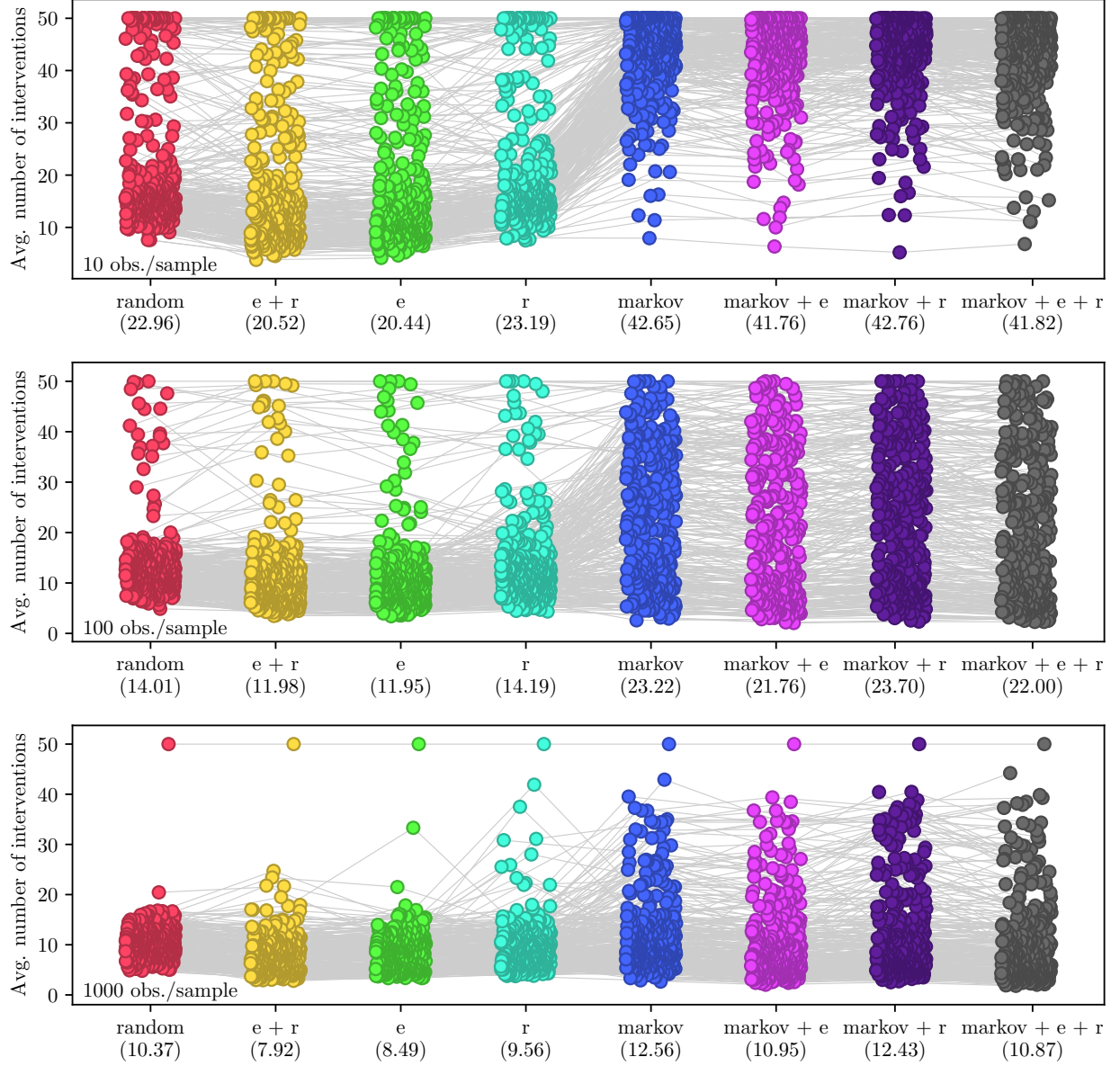


Figure 5: (finite regime) Average number of interventions until the causal predictors are identified, for each one of the 300 SEMs. Each SEM is represented by a dot and connected across policies by a grey line. The total average of interventions employed by each policy is given below its label. The policies are run for a maximum of 50 iterations, so for SEMs at this number of iterations the policies likely did not yield the correct estimate. For the policies that employ the Markov blanket estimate, we can clearly see the effect that the sample size has on the quality of such estimate, and by extension, on the policies' performance.

## 6 Discussion

In Proposition 2 we show that a direct cause appears on at least half of all stable sets. We use this result to construct an intervention selection policy for invariant causal prediction, that is consistent in the number of interventions and outperforms the random policy in both population and finite-regime experiments.

Many interesting questions remain. ICP does not require knowledge of the intervention locations in each environment, which makes it robust to interventions with off-target effects i.e. effects on variables other than the target. On the other hand, one might ask if, since we know the intervention location, we are throwing away useful information. Of the proposed policies, only the ones that use the *empty-set* heuristic use this knowledge, but do so at the cost of becoming inconsistent. Furthermore, this issue of inconsistency arises from a binary, deterministic and irreversible decision, and as such can be framed as a Noisy Generalized Binary Search problem [Nowak 2009], for which a body of research already exists. Finally, the result from Proposition 2 is quite general in the sense that it makes no assumptions on the function class or noise distributions of the SEM. As such, it is interesting to ask if it could also be used to inform intervention selection policies for more general extensions of invariant causal prediction, such as non-linear ICP [Heinze-Deml, Peters, and Meinshausen 2018] or ICP for sequential data [Pfister, Bühlmann, and Peters 2019].

## 7 Acknowledgments

I thank Christina Heinze-Deml, Niklas Pfister and Jonas Peters for the valuable discussion and comments on the manuscript.

The research leading to these results was supported by a grant from the "la Caixa" Foundation (ID 100010434), with code LCF/BQ/EU18/11650051.

## References

- [1] Raj Agrawal, Chandler Squires, Karren Yang, Karthikeyan Shanmugam, and Caroline Uhler. "ABCD-Strategy: Budgeted Experimental Design for Targeted Causal Structure Discovery". In: *The 22nd International Conference on Artificial Intelligence and Statistics*. 2019, pp. 3400–3409.
- [2] Steen A Andersson, David Madigan, Michael D Perlman, et al. "A characterization of Markov equivalence classes for acyclic digraphs". In: *The Annals of Statistics* 25.2 (1997), pp. 505–541.
- [3] Frederick Eberhardt. "Almost optimal intervention sets for causal discovery". In: *Proceedings of the Twenty-Fourth Conference on Uncertainty in Artificial Intelligence*. 2008, pp. 161–168.
- [4] AmirEmad Ghassami, Saber Salehkaleybar, Negar Kiyavash, and Elias Bareinboim. "Budgeted Experiment Design for Causal Structure Learning". In: *International Conference on Machine Learning*. 2018, pp. 1719–1728.
- [5] Alain Hauser and Peter Bühlmann. "Two optimal strategies for active learning of causal models from interventional data". In: *International Journal of Approximate Reasoning* 55.4 (2014), pp. 926–939.
- [6] Yang-Bo He and Zhi Geng. "Active learning of causal networks with intervention experiments and optimal designs". In: *Journal of Machine Learning Research* 9.Nov (2008), pp. 2523–2547.
- [7] Christina Heinze-Deml, Jonas Peters, and Nicolai Meinshausen. "Invariant causal prediction for non-linear models". In: *Journal of Causal Inference* 6.2 (2018).
- [8] Andrés R Masegosa and Serafin Moral. "An interactive approach for Bayesian network learning using domain/expert knowledge". In: *International Journal of Approximate Reasoning* 54.8 (2013), pp. 1168–1181.
- [9] Kevin P Murphy. "Active learning of causal Bayes net structure". In: (2001).
- [10] Robert Osazuwa Ness, Karen Sachs, Parag Mallick, and Olga Vitek. "A Bayesian active learning experimental design for inferring signaling networks". In: *International Conference on Research in Computational Molecular Biology*. Springer. 2017, pp. 134–156.

- [11] Robert Nowak. “Noisy generalized binary search”. In: *Advances in neural information processing systems*. 2009, pp. 1366–1374.
- [12] Judea Pearl. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann, 1988.
- [13] Judea Pearl et al. “Causal inference in statistics: An overview”. In: *Statistics surveys* 3 (2009), pp. 96–146.
- [14] J Peters, J Mooij, D Janzing, and B Schölkopf. “Identifiability of causal graphs using functional models”. In: *27th Conference on Uncertainty in Artificial Intelligence (UAI 2011)*. AUAI Press. 2011, pp. 589–598.
- [15] Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. “Causal inference by using invariant prediction: identification and confidence intervals”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 78.5 (2016), pp. 947–1012.
- [16] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations and learning algorithms*. MIT press, 2017.
- [17] Niklas Pfister, Peter Bühlmann, and Jonas Peters. “Invariant causal prediction for sequential data”. In: *Journal of the American Statistical Association* 114.527 (2019), pp. 1264–1276.
- [18] Niklas Pfister, Evan G. Williams, Jonas Peters, Ruedi Aebersold, and Peter Bühlmann. *Stabilizing Variable Selection and Regression*. 2019. arXiv: 1911.01850 [stat.ME].
- [19] Burr Settles. *Active learning literature survey*. Tech. rep. University of Wisconsin-Madison Department of Computer Sciences, 2009.
- [20] Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. *Causation, prediction, and search*. 2000.
- [21] Robert Tibshirani. “Regression shrinkage and selection via the lasso”. In: *Journal of the Royal Statistical Society: Series B (Methodological)* 58.1 (1996), pp. 267–288.
- [22] Simon Tong and Daphne Koller. “Active learning for structure in Bayesian networks”. In: *International joint conference on artificial intelligence*. Vol. 17. 1. Citeseer. 2001, pp. 863–869.
- [23] Thomas Verma and Judea Pearl. “Equivalence and synthesis of causal models”. In: *Proceedings of the Sixth Annual Conference on Uncertainty in Artificial Intelligence*. Elsevier Science Inc. 1990, pp. 255–270.
- [24] Stephen G Walker. “Bayesian inference with misspecified models”. In: *Journal of Statistical Planning and Inference* 143.10 (2013), pp. 1621–1633.

## A Proofs for section 2

The following two lemmas are used in the proofs of later statements.

**Lemma 5.** *Let  $\mathcal{E}_t, \mathcal{E}_{t+1}$  be sets of observed environments such that  $\mathcal{E}_t \subseteq \mathcal{E}_{t+1}$ . Then,*

$$N_{t+1}^{\text{int}} = N_t^{\text{int}} \setminus \text{DE}(CH_{t+1}^{\text{int}}(Y) \setminus CH_t^{\text{int}}(Y))$$

i.e.  $N_{t+1}^{\text{int}}$  is  $N_t^{\text{int}}$  without the descendants of the newly intervened children of  $Y$ .

*Proof.*

$$N_{t+1}^{\text{int}} := \{1, \dots, p\} \setminus \text{DE}(CH_{t+1}^{\text{int}}(Y)).$$

Since  $CH_t^{\text{int}}(Y) \subseteq CH_{t+1}^{\text{int}}(Y)$ , we can write,

$$\begin{aligned} N_{t+1}^{\text{int}} &= \{1, \dots, p\} \setminus \text{DE}(CH_t^{\text{int}}(Y) \cup (CH_{t+1}^{\text{int}}(Y) \setminus CH_t^{\text{int}}(Y))) \\ &= \{1, \dots, p\} \setminus \text{DE}(CH_t^{\text{int}}(Y)) \setminus \text{DE}(CH_{t+1}^{\text{int}}(Y) \setminus CH_t^{\text{int}}(Y)) \\ &= N_t^{\text{int}} \setminus \text{DE}(CH_{t+1}^{\text{int}}(Y) \setminus CH_t^{\text{int}}(Y)). \end{aligned}$$

□

Note that Lemma 5 implies that  $N_t^{\text{int}} = N_{t+1}^{\text{int}} \cup \text{DE}(CH_{t+1}^{\text{int}}(Y) \setminus CH_t^{\text{int}}(Y))$ .

**Lemma 6** (the stable blanket shrinks with informative environments). *Let  $\mathcal{E}_t, \mathcal{E}_{t+1}$  be sets of observed environments such that  $\mathcal{E}_t \subseteq \mathcal{E}_{t+1}$ . Then,*

$$SB_{t+1}(Y) \subseteq SB_t(Y).$$

*Proof.* From Theorem 1 we have that

$$\begin{aligned} SB_t(Y) &= \text{PA}(Y) \cup (\text{CH}(Y) \cap N_t^{\text{int}}) \cup \text{PA}(\text{CH}(Y) \cap N_t^{\text{int}}), \\ SB_{t+1}(Y) &= \text{PA}(Y) \cup (\text{CH}(Y) \cap N_{t+1}^{\text{int}}) \cup \text{PA}(\text{CH}(Y) \cap N_{t+1}^{\text{int}}). \end{aligned}$$

By lemma 1 we have that  $N_{t+1}^{\text{int}} \subseteq N_t^{\text{int}}$ , and it follows

$$\begin{aligned} \text{CH}(Y) \cap N_{t+1}^{\text{int}} &\subseteq \text{CH}(Y) \cap N_t^{\text{int}}, \\ \text{PA}(\text{CH}(Y) \cap N_{t+1}^{\text{int}}) &\subseteq \text{PA}(\text{CH}(Y) \cap N_t^{\text{int}}) \end{aligned}$$

i.e.  $SB_{t+1}(Y) \subseteq SB_t(Y)$ .

□

**Proposition 1** (effect of interventions on the stable blanket). *Let  $\mathcal{E}_t, \mathcal{E}_{t+1}$  be sets of observed environments such that  $\mathcal{E}_t \subseteq \mathcal{E}_{t+1}$ . Let  $I_{t+1}$  be the variables intervened in  $e_{t+1}$ , and let  $S = CH_{t+1}^{\text{int}}(Y) \setminus CH_t^{\text{int}}(Y) = I_{t+1} \cap CH(Y) \cap N_t^{\text{int}}$  be the children of the response intervened on for the first time in  $e_{t+1}$ . Then,*

$$\begin{aligned} SB_t(Y) &= SB_{t+1}(Y) \cup \\ &\quad (\text{DE}(S) \cap SB_t(Y)) \cup \\ &\quad \{j \in SB_t(Y) : j \notin \text{DE}(S) \wedge j \notin (\text{CH}(Y) \cap N_t^{\text{int}}) \wedge \text{CH}(j) \subseteq \text{DE}(S)\}. \end{aligned}$$

*Proof.* " $\subseteq$ ". Let  $j \in SB_t(Y)$ , then either

- a)  $j \in \text{PA}(Y) \implies j \in SB_{t+1}(Y)$ .
- b)  $j \in \text{CH}(Y) \cap N_t^{\text{int}}$ , and by Lemma 5 we have that

$$\begin{aligned} j \in \text{CH}(Y) \cap (N_{t+1}^{\text{int}} \cup \text{DE}(S)) &\implies j \in \text{CH}(Y) \cap N_{t+1}^{\text{int}} \vee j \in \text{CH}(Y) \cap \text{DE}(S) \\ &\implies j \in SB_{t+1}(Y) \vee j \in \text{DE}(S) \cap SB_t(Y). \end{aligned}$$



- c)  $j \notin PA(Y) \wedge j \notin CH(Y) \cap N_t^{\text{int}}$ , which means that  $CH(j) \subseteq CH(Y) \cap N_t^{\text{int}}$ . If additionally,  $j \in DE(S) \implies j \in DE(S) \cap SB_t(Y)$ .

Otherwise, by Lemma 5 we have that

$$CH(j) \subseteq CH(Y) \cap (N_{t+1}^{\text{int}} \cup DE(S)),$$

which means that either

$$\exists i \in CH(j) : i \in CH(Y) \cap N_{t+1}^{\text{int}},$$

which implies that  $j \in SB_{t+1}(Y)$ , or

$$i \in DE(S) \quad \forall i \in CH(j),$$

by which

$$j \in \{j \in SB_t(Y) : j \notin DE(S) \wedge j \notin (CH(Y) \cap N_t^{\text{int}}) \wedge CH(j) \subseteq DE(S)\}.$$

" $\supseteq$ " By Lemma 6 we have that  $SB_{t+1}(Y) \subseteq SB_t(Y)$ , and the other two sets of the right hand side are subsets of  $SB_t(Y)$  by definition.  $\square$

**Lemma 1** (interventions on upstream variables). *Let  $\mathcal{E}$  be a set of observed environments and let  $j \in AN(Y)$  such that it is directly intervened in  $\mathcal{E}$ . Then,*

$$S \subseteq \{1, \dots, p\} \text{ is intervention stable} \implies S \cap DE(j) \cap AN(Y) \neq \emptyset.$$

*Proof.* Assume  $S \subseteq \{1, \dots, p\}$  is an intervention stable set such that  $S \cap DE(j) \cap AN(Y) = \emptyset$ , and let  $I^j$  denote the direct intervention on  $j$ . Then, there is a path  $I^j \rightarrow j \rightarrow \dots \rightarrow Y$ , which is unblocked by  $S$ .  $\square$

**Lemma 2** (intervened parents appear on all intervention stable sets). *Let  $\mathcal{E}$  be a set of observed environments and let  $j \in PA(Y)$  such that it is directly intervened in  $\mathcal{E}$ . Then,*

$$S \subseteq \{1, \dots, p\} \text{ is intervention stable} \implies j \in S.$$

*Proof.* Assume  $S \subseteq \{1, \dots, p\}$  is an intervention stable set such that  $j \notin S$ , and let  $I^j$  denote the direct intervention on  $j$ . Then, there is a path  $I^j \rightarrow j \rightarrow Y$  which is unblocked by  $S$ .  $\square$

**Lemma 3** (sets containing descendants of directly intervened children are unstable). *Let  $i \in CH^{\text{int}}(Y)$ . Then, any set  $S \subseteq \{1, \dots, d\}$  which contains descendants of  $i$  is not intervention stable.*

*Proof.* Let  $I^i$  denote the direct intervention on  $i$ , and let  $S \subseteq \{1, \dots, p\} : S \cap DE(i) \neq \emptyset$ . Then, the path  $Y \rightarrow i \leftarrow I^i$  is not blocked by  $S$ .  $\square$

**Lemma 4** (stability of parents of intervened children). *Let  $\mathcal{E}_t, \mathcal{E}_{t+1}$  be sets of observed environments such that  $\mathcal{E}_t \subseteq \mathcal{E}_{t+1}$ . Assume an intervention occurred such that  $j \in SB_t(Y)$  but  $j \notin SB_{t+1}(Y)$ , and let  $S = CH_{t+1}^{\text{int}}(Y) \setminus CH_t^{\text{int}}(Y)$  be the set of newly intervened children. Then, if  $SB_{t+1}(Y) \cup \{j\}$  is intervention stable,*

$$j \in \{j \in SB_t(Y) : j \notin DE(S) \wedge j \notin (CH(Y) \cap N_t^{\text{int}}) \wedge CH(j) \subseteq DE(S)\}.$$

*Proof.* Since  $j \in SB_t(Y)$  but  $j \notin SB_{t+1}(Y)$ , by Proposition 1 we have that either

a)  $j \in DE(S) \cap SB_t(Y)$ , which contradicts Lemma 3, or

b)  $j \in \{j \in SB_t(Y) : j \notin DE(S) \wedge j \notin (CH(Y) \cap N_t^{\text{int}}) \wedge CH(j) \subseteq DE(S)\},$

i.e.  $j$  is a parent of descendants of intervened children.  $\square$

**Proposition 2** (parents appear on at least half of all stable sets). *Let  $\mathcal{E}$  be any set of observed environments. Then, for any  $j \in \{1, \dots, p\}$ ,*

$$r_{\mathcal{E}}(j) < 1/2 \implies j \notin PA(Y).$$

*Proof.* We will prove the equivalent statement  $j \in \text{PA}(Y) \implies r_{\mathcal{E}}(j) \geq 1/2$ . For any  $i \in \{1, \dots, p\}$  we have that

$$r_{\mathcal{E}}(i) = \frac{|\{S \in \mathbb{S}_{\mathcal{E}} : i \in S\}|}{|\{S \in \mathbb{S}_{\mathcal{E}} : i \in S\}| + |\{S \in \mathbb{S}_{\mathcal{E}} : i \notin S\}|},$$

and therefore

$$r_{\mathcal{E}}(i) \geq 1/2 \iff |\{S \in \mathbb{S}_{\mathcal{E}} : i \in S\}| \geq |\{S \in \mathbb{S}_{\mathcal{E}} : i \notin S\}|. \quad (3)$$

We will show that for any  $j \in \text{PA}(Y)$ , and any intervention stable set  $S$  such that  $j \notin S$ , the set  $S \cup \{j\}$  is also intervention stable, satisfying the right hand side of Equation 3. To do this, we will use the fact that  $S$  d-separates the response from all interventions, and show that the same is true for  $S \cup \{j\}$ , making it intervention stable.

Let  $I$  denote an intervention on a variable  $i$ . For every path connecting  $Y$  and the intervention, either

- (i)  $j$  appears in the path as a collider,
- (ii)  $j$  appears in the path but not as a collider,
- (iii)  $j$  does not appear in the path but is downstream of a collider, or
- (iv)  $j$  does not appear in the path and is not downstream of a collider.

If  $S$  blocks paths of type (ii) and (iv),  $S \cup \{j\}$  also does. Assume now there is a path of type (i) or (iii) which is blocked under  $S$  but active under  $S \cup \{j\}$ . This implies that such path is blocked by a collider  $c$  such that  $j \in \text{DE}(c)$ ,  $S \cap \text{DE}(c) = \emptyset$ , and there exists a path  $Y \leftarrow j \leftarrow \dots \leftarrow c \leftarrow \dots \leftarrow i \leftarrow I$  which is active under  $S$ , i.e.  $S \notin \mathbb{S}_{\mathcal{E}}$ .

Therefore, for all  $S \in \mathbb{S}_{\mathcal{E}}$  such that  $j \notin S$ , we have that  $S \cup \{j\} \in \mathbb{S}_{\mathcal{E}}$ , and

$$|\{S \in \mathbb{S}_{\mathcal{E}} : j \in S\}| \geq |\{S \in \mathbb{S}_{\mathcal{E}} : j \notin S\}| \implies r_{\mathcal{E}}(j) \geq 1/2.$$

□

## B Proofs for section 3

**Proposition 3** (intervention stable sets are plausible causal predictors). *Let  $\mathcal{E}$  be a set of observed environments. Then, for all intervention stable sets  $S \subseteq \{1, \dots, p\}$ , it holds that  $S \in \mathbb{C}_{\mathcal{E}}$ .*

*Proof.* The following is based on proof of proposition 3 in [Pfister et al. 2019].

Let  $\mathcal{E}$  a set of observed environments, and let  $S \in \mathbb{S}_{\mathcal{E}}$  be an intervention stable set. From [Jonas Peters, Bühlmann, and Meinshausen 2016] we know that  $S$  is a set of plausible causal predictors iff  $Y^e | X_S^e$  remains invariant for all environments  $e \in \mathcal{E}$ . Starting from setting 1, introduce an auxiliary random variable  $E$  taking values in  $\mathcal{E}$  with equal probability (for simplicity). To model the environments we construct an extended SEM  $\mathcal{S}_{\text{full}}^{\mathcal{E}}$ , where the variable  $E$  appears on the assignments of the intervention variables  $I$ , and the assignments of the remaining variables remain as in  $\mathcal{S}^{\mathcal{E}}$ . As such, in  $\mathcal{G}(\mathcal{S}_{\text{full}}^{\mathcal{E}})$   $E$  is a source node with only edges into the variables in  $I$ . The SEM  $\mathcal{S}_{\text{full}}^{\mathcal{E}}$  induces a distribution  $P_{\text{full}}$  over  $(E, I, X, Y)$ , which under setting 1 has a density  $p$  that factorizes with respect to a product measure. Furthermore, since  $P_{\text{full}}$  satisfies the Markov properties [Pearl et al. 2009] and  $S$  d-separates the response from all the intervention variables in  $I$ , it holds that  $E \perp\!\!\!\perp Y \mid X_S \sim P_{\text{full}}$ . Therefore, for every environment  $e \in \mathcal{E}$ , we have that

$$\begin{aligned} p(Y^e = y \mid X_S^e = x) &= p(Y = y \mid X_S = x, E = e) \\ &= \frac{p(Y = y, E = e \mid X_S = x)}{p(E = e \mid X_S = x)} \\ &= \frac{p(Y = y \mid X_S = x)p(E = e \mid X_S = x)}{p(E = e \mid X_S = x)} \\ &= p(Y = y \mid X_S = x), \end{aligned}$$

and  $Y^e \mid X_S^e$  remains invariant for all environments  $e \in \mathcal{E}$ . □

# Declaration of Originality

The signed declaration of originality is a component of every semester paper, Bachelor's thesis, Master's thesis and any other degree paper undertaken during the course of studies, including the respective electronic versions.

Lecturers may also require a declaration of originality for other written papers compiled for their courses.

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the supervisor .

**Title of work** (in block letters):

TOWARDS ACTIVE ICP: EXPERIMENT SELECTION THROUGH STABILITY

**Authored by** (in block letters):

*For papers written by groups the names of all authors are required.*

**Name(s):**

GIANELLA MARTIN

**First name(s):**

JUAN LUIS

With my signature I confirm that

- I have committed none of the forms of plagiarism described in the **Citation etiquette** information sheet.
- I have documented all methods, data and processes truthfully.
- I have not manipulated any data.
- I have mentioned all persons who were significant facilitators of the work .
- I am aware that the work may be screened electronically for plagiarism.
- I have understood and followed the guidelines in the document *Scientific Works in Mathematics*.

**Place, date:**

ZÜRICH, 07/02/2020

**Signature(s):**

*For papers written by groups the names of all authors are required. Their signatures collectively guarantee the entire content of the written paper.*