

LEARNING COUNTERFACTUALLY INVARIANT PREDICTORS

Anonymous authors

Paper under double-blind review

ABSTRACT

Notions of counterfactual invariance (CI) have proven essential for predictors that are fair, robust, and generalizable in the real world. We propose graphical criteria that yield a sufficient condition for a predictor to be counterfactually invariant in terms of a conditional independence in the observational distribution. In order to learn such predictors, we propose a model-agnostic framework, called Counterfactually Invariant Prediction (CIP), building on the Hilbert-Schmidt Conditional Independence Criterion (HSCIC), a kernel-based conditional dependence measure. Our experimental results demonstrate the effectiveness of CIP in enforcing counterfactual invariance across various simulated and real-world datasets including scalar and multi-variate settings.

1 INTRODUCTION

Invariance, or equivariance to certain data transformations, has proven essential in numerous applications of machine learning (ML), since it can lead to better generalization capabilities (Arjovsky et al., 2019; Chen et al., 2020; Bloem-Reddy & Teh, 2020). For instance, in image recognition, predictions ought to remain unchanged under scaling, translation, or rotation of the input image. Data augmentation, an early heuristic to promote such invariances, has become indispensable for successfully training deep neural networks (DNNs) (Shorten & Khoshgoftaar, 2019; Xie et al., 2020). Well-known examples of “invariance by design” include convolutional neural networks (CNNs) for translation invariance (Krizhevsky et al., 2012), group equivariant NNs for general group transformations (Cohen & Welling, 2016), recurrent neural networks (RNNs) and transformers for sequential data (Vaswani et al., 2017), DeepSet (Zaheer et al., 2017) for sets, and graph neural networks (GNNs) for different types of geometric structures (Battaglia et al., 2018).

Many applications in modern ML, however, call for arguably stronger notions of invariance based on causality. This case has been made for image classification, algorithmic fairness (Hardt et al., 2016; Mitchell et al., 2021), robustness (Bühlmann, 2020), and out-of-distribution generalization (Lu et al., 2021). The goal is invariance with respect to hypothetical manipulations of the data generating process (DGP). Various works develop methods that assume observational distributions (across environments or between training and test) to be governed by shared causal mechanisms, but differ due to various types of distribution shifts encoded by the causal model (Peters et al., 2016; Heinze-Deml et al., 2018; Rojas-Carulla et al., 2018; Arjovsky et al., 2019; Bühlmann, 2020; Subbaswamy et al., 2022; Yi et al., 2022; Makar et al., 2022). Typical goals include to train predictors invariant to such shifts, to learn about causal mechanisms and to improve robustness against spurious correlations or out of distribution generalization. The term “counterfactual invariance” has also been used in other out of distribution learning contexts unrelated to our task, e.g., to denote invariance to certain symmetry transformations (Mouli & Ribeiro, 2022).

While we share the broader motivation, these works are orthogonal to ours, because even though counterfactual distributions are also generated from the same causal model, they are fundamentally different from such shifts (Peters et al., 2017). Intuitively, counterfactuals are about events that did not, but could have happened had circumstances been different in a controlled way. A formal discussion of what we mean by counterfactuals is required to properly position our work in the existing literature and describe our contributions.

2 PROBLEM SETTING AND RELATED WORK

2.1 PRELIMINARIES AND TERMINOLOGY

Definition 2.1 (Structural causal model (SCM)). A structural causal model is a tuple $\mathcal{S} = (\mathbf{U}, \mathbf{V}, F, \mathbb{P}_{\mathbf{U}})$ such that \mathbf{U} is a set of background variables that are exogenous to the model; \mathbf{V}

is a set of observable (endogenous) variables; $F = \{f_V\}_{V \in \mathbf{V}}$ is a set of functions from (the domains of) $\text{pa}(V) \cup U_V$ to (the domain of) V , where $U_V \subset \mathbf{U}$ and $\text{pa}(V) \subseteq \mathbf{V} \setminus \{V\}$ such that $V = f_V(\text{pa}(V), U_V)$; $\mathbb{P}_{\mathbf{U}}$ is a probability distribution over the domain of \mathbf{U} . Further, the subsets $\text{pa}(V) \subseteq \mathbf{V} \setminus \{V\}$ are chosen such that the graph \mathcal{G} over \mathbf{V} where the edge $V' \rightarrow V$ is in \mathcal{G} if and only if $V' \in \text{pa}(V)$ is a directed acyclic graph (DAG).

Observational distribution. An SCM implies a unique observational distribution over \mathbf{V} , which can be thought of as being generated by transforming the distribution over $\mathbb{P}_{\mathbf{U}}$ via the deterministic functions in F iteratively to obtain a distribution over \mathbf{V} .¹

Interventions. Given a variable $A \in \mathbf{V}$, an intervention $A \leftarrow a$ amounts to replacing f_A in F with the constant function setting A to a . This yields a new SCM, which induces the *interventional distribution* under intervention $A \leftarrow a$.² Similarly, we can intervene on multiple variables $\mathbf{V} \supseteq \mathbf{A} \leftarrow \mathbf{a}$. For an outcome (or prediction target) variable $\mathbf{Y} \subset \mathbf{V}$, we then write $\mathbf{Y}_{\mathbf{a}}^*$ for the outcome in the intervened SCM, also called *potential outcome*. Note that the interventional distribution $\mathbb{P}_{\mathbf{Y}_{\mathbf{a}}^*}(\mathbf{y})$ differs in general from the conditional distribution $\mathbb{P}_{\mathbf{Y}|\mathbf{A}}(\mathbf{y} | \mathbf{a})$.³ This is typically the case when \mathbf{Y} and \mathbf{A} have a shared ancestor, i.e., they are confounded. In interventional distributions, potential outcomes are random variables via the exogenous variables \mathbf{u} , i.e., $\mathbf{Y}_{\mathbf{a}}^*(\mathbf{u})$ where $\mathbf{u} \sim \mathbb{P}_{\mathbf{U}}$. Hence, interventions capture “population level” properties, i.e., the action is performed for all units \mathbf{u} .

Counterfactuals. Counterfactuals capture what happens under interventions for a “subset” of possible units \mathbf{u} that are compatible with observations $\mathbf{W} = \mathbf{w}$ for a subset of observed variables $\mathbf{W} \subseteq \mathbf{V}$. This can be described in a three step procedure. (i) *Abduction*: We restrict our attention to units compatible with the observations, i.e., consider the new SCM $\mathcal{S}^{\mathbf{w}} = (\mathbf{U}, \mathbf{V}, F, \mathbb{P}_{\mathbf{U}|\mathbf{W}=\mathbf{w}})$. (ii) *Intervention*: Within $\mathcal{S}^{\mathbf{w}}$, perform an intervention $\mathbf{A} \leftarrow \mathbf{a}$ on some variables \mathbf{A} (which need not be disjoint from \mathbf{W}). (iii) *Prediction*: Finally, we are typically interested in the outcome \mathbf{Y} in an interventional distribution of $\mathcal{S}^{\mathbf{w}}$, which we denote by $\mathbb{P}_{\mathbf{Y}_{\mathbf{a}}^*|\mathbf{W}=\mathbf{w}}(\mathbf{y})$ and call a *counterfactual distribution*: “Given that we have observed $\mathbf{W} = \mathbf{w}$, what would \mathbf{Y} have been had we set $\mathbf{A} \leftarrow \mathbf{a}$, instead of the value \mathbf{A} has actually taken?” Counterfactuals capture properties of a “subpopulation” $\mathbf{u} \sim \mathbb{P}_{\mathbf{U}|\mathbf{W}=\mathbf{w}}$ compatible with the observations.⁴ Even for granular \mathbf{W} , there may be multiple units \mathbf{u} in the support of this distribution. In contrast, “unit level counterfactuals” often considered in philosophy contrast $\mathbf{Y}_{\mathbf{a}}^*(\mathbf{u})$ with $\mathbf{Y}_{\mathbf{a}'}^*(\mathbf{u})$ for a single unit \mathbf{u} . Such unit level counterfactuals are too fine-grained in our setting. Hence, our used definition of counterfactual invariance is:

Definition 2.2 (Counterfactual invariance). Let \mathbf{A}, \mathbf{W} be (not necessarily disjoint) sets of nodes in a given SCM. Then, \mathbf{Y} is *counterfactually invariant in \mathbf{A} w.r.t. \mathbf{W}* if $\mathbb{P}_{\mathbf{Y}_{\mathbf{a}}^*|\mathbf{W}=\mathbf{w}}(\mathbf{y}) = \mathbb{P}_{\mathbf{Y}_{\mathbf{a}'}^*|\mathbf{W}=\mathbf{w}}(\mathbf{y})$ almost surely, for all \mathbf{a}, \mathbf{a}' in the domain of \mathbf{A} and all \mathbf{w} in the domain of \mathbf{W} .⁵

Predictors in SCMs. Ultimately, we aim at learning a predictor $\hat{\mathbf{Y}}$ for the outcome \mathbf{Y} . Originally, the predictor $\hat{\mathbf{Y}}$ is not part of the DGP, because we get to learn $f_{\hat{\mathbf{Y}}}$ from data. Using supervised learning, the predictor $f_{\hat{\mathbf{Y}}}$ depends both on the chosen inputs $\mathbf{X} \subset \mathbf{V}$ as well as the target \mathbf{Y} . However, once $f_{\hat{\mathbf{Y}}}$ is fixed, it is a deterministic function with arguments $\mathbf{X} \subset \mathbf{V}$, so $(\mathbf{U}, \mathbf{V} \cup \{\hat{\mathbf{Y}}\}, F \cup \{f_{\hat{\mathbf{Y}}}\}, \mathbb{P}_{\mathbf{U}})$ is a valid SCM and we can consider $\hat{\mathbf{Y}}$ an observed variable with incoming arrows from only \mathbf{X} . Hence, the definition of counterfactual invariance can be applied to the predictor $\hat{\mathbf{Y}}$.

Kernel mean embeddings (KME). Our method relies on kernel mean embeddings (KMEs). We describe the main concepts pertaining KMEs and refer the reader to [Smola et al. \(2007\)](#); [Schölkopf et al. \(2002\)](#); [Berlinet & Thomas-Agnan \(2011\)](#); [Muandet et al. \(2017\)](#) for details. Fix a measurable space \mathcal{Y} with respect to a σ -algebra $\mathcal{F}_{\mathcal{Y}}$, and consider a probability measure \mathbb{P} on the space $(\mathcal{Y}, \mathcal{F}_{\mathcal{Y}})$. Let \mathcal{H} be a reproducing kernel Hilbert space (RKHS) with a bounded kernel $k_{\mathbf{Y}}: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$, i.e., $k_{\mathbf{Y}}$ is such that $\sup_{\mathbf{y} \in \mathcal{Y}} k(\mathbf{y}, \mathbf{y}) < \infty$. The kernel mean embedding $\mu_{\mathbb{P}}$ of \mathbb{P} is defined as the expected value of the function $k(\cdot, \mathbf{y})$ with respect to \mathbf{y} , i.e., $\mu_{\mathbb{P}} := \mathbb{E}[k(\cdot, \mathbf{y})]$. The definition

¹Note that all randomness stems from $\mathbb{P}_{\mathbf{U}}$. The observational distribution is well-defined and unique, essentially because every DAG allows for a topological order.

²The observational distribution in an intervened SCM is called interventional distribution of the base SCM.

³We use \mathbb{P} for distributions (common in the kernel literature) and $\mathbf{Y}_{\mathbf{a}}^*$ instead of the do notation.

⁴Note that conditioning in an interventional distribution is different from a counterfactual and our notation is quite subtle here $\mathbb{P}_{\mathbf{Y}_{\mathbf{a}}^*}(\mathbf{y} | \mathbf{W} = \mathbf{w}) \neq \mathbb{P}_{\mathbf{Y}_{\mathbf{a}'}^*|\mathbf{W}=\mathbf{w}}(\mathbf{y})$.

⁵With an abuse of notation, if $\mathbf{W} = \emptyset$ then the requirement of conditional counterfactual invariance becomes $\mathbb{P}_{\mathbf{Y}_{\mathbf{a}}^*}(\mathbf{y}) = \mathbb{P}_{\mathbf{Y}_{\mathbf{a}'}^*}(\mathbf{y})$ almost surely, for all \mathbf{a}, \mathbf{a}' in the domain of \mathbf{A} .

of KMEs can be extended to conditional distributions (Fukumizu et al., 2013; Grünewälder et al., 2012; Song et al., 2009; 2013). Consider two random variables \mathbf{Y} , \mathbf{S} , and denote with $(\Omega_{\mathbf{Y}}, \mathcal{F}_{\mathbf{Y}})$ and $(\Omega_{\mathbf{S}}, \mathcal{F}_{\mathbf{S}})$ the respective measurable spaces. These random variables induce a probability measure $\mathbb{P}_{\mathbf{Y}, \mathbf{S}}$ in the product space $\Omega_{\mathbf{Y}} \times \Omega_{\mathbf{S}}$. Let $\mathcal{H}_{\mathbf{Y}}$ be a RKHS with a bounded kernel $k_{\mathbf{Y}}(\cdot, \cdot)$ on $\Omega_{\mathbf{Y}}$. We define the KME of a conditional distribution $\mathbb{P}_{\mathbf{Y}|\mathbf{S}}(\cdot | \mathbf{s})$ via $\mu_{\mathbf{Y}|\mathbf{S}=\mathbf{s}} := \mathbb{E}[k_{\mathbf{Y}}(\cdot, \mathbf{y}) | \mathbf{S} = \mathbf{s}]$. Here, the expected value is taken over \mathbf{y} . KMEs of conditional measures can be estimated from samples (Grünewälder et al., 2012). Pogodin et al. (2022) recently proposed an efficient kernel-based regularizer for learning features of input data that allow for estimating a target while being conditionally independent of a distractor given the target. Since CIP ultimately enforces conditional independence (see Theorem 3.2), we believe it could further benefit from leveraging the efficiency and convergence properties of their technique, which we leave for future work.

2.2 RELATED WORK AND CONTRIBUTIONS

While we focus on counterfactuals in the SCM framework (Pearl, 2000; Peters et al., 2016), there are different incompatible frameworks to describe counterfactuals (von Kügelgen et al., 2022; Dorr, 2016; Woodward, 2021), which may give rise to orthogonal notions of counterfactual invariance.

Research on algorithmic fairness has explored a plethora of causal “invariance” notions with the goal of achieving fair predictors (Loftus et al., 2018; Carey & Wu, 2022; Plecko & Bareinboim, 2022). Kilbertus et al. (2017) conceptually introduce a notion based on group-level interventions, which has been refined to take into account more granular context by Salimi et al. (2019); Galhotra et al. (2022), who then obtain fair predictors by viewing it as a database repair problem or a causal feature selection problem, respectively. A counterfactual-level definition was proposed by Kusner et al. (2017) and followed up by path-specific counterfactual notions (Nabi & Shpitser, 2018; Chiappa, 2019), where the protected attribute may take different values along different paths to the outcome. Recently, Dutta et al. (2021) developed an information theoretic framework to decompose the overall causal influence allowing for exempted variables and properly dealing with synergies across different paths.

Our focus is on counterfactuals because they are fundamentally more expressive than mere interventions (Pearl, 2000; Bareinboim et al., 2022), but do not require a fine-grained path- or variable-level judgment of “allowed” and “disallowed” paths or variables, which may be challenging to devise in practice. Since CI already requires strong assumptions, we leave path-specific counterfactuals—even more challenging in terms of identifiability (Avin et al., 2005)—for future work. While our Definition 2.2 requires equality in distribution, Veitch et al. (2021) suggest a definition of a counterfactually invariant predictor $f_{\hat{\mathbf{Y}}}$ which requires almost sure equality of $\hat{\mathbf{Y}}_{\mathbf{a}}^*$ and $\hat{\mathbf{Y}}_{\mathbf{a}'}^*$, where we view $\hat{\mathbf{Y}}$ as an observed variable in the SCM as described above. Fawkes & Evans (2023) recently shed light on the connection between almost sure CI (a.s.-CI), distributional CI (\mathcal{D} -CI) as in Definition 2.2, and CI of predictors (\mathcal{F} -CI) showing that $f_{\hat{\mathbf{Y}}}$ being \mathcal{F} -CI is equivalent to $\hat{\mathbf{Y}}$ being \mathcal{D} -CI conditioned on \mathbf{X} , rendering it also equivalent to counterfactual fairness (Kusner et al., 2017).

Inspired by problems in natural language processing (NLP), Veitch et al. (2021) aim at “stress-testing” models for spurious correlations. It differs from our work in that they (i) focus only on two specific graphs, and (ii) provide a *necessary* but not sufficient criterion for CI in terms of a conditional independence. Their method enforces the conditional independence via maximum mean discrepancy (MMD) (in *discrete settings only*). However, enforcing a consequence of CI, does not necessarily improve CI. Indeed, Fawkes & Evans (2023, Prop. 4.4) show that while a.s.-CI implies certain conditional independencies, no set of conditional independencies implies any bounds on the difference in counterfactuals. On the contrary, the weaker notion of \mathcal{D} -CI in Definition 2.2 can in fact be equivalent to conditional independencies in the observational distribution (Fawkes & Evans, 2023, Lem. A.3). Albeit only proved for special cases where the counterfactual distribution is identifiable, this opens the possibility for sufficient graphical criteria for distributional CI in any graph.

Contributions. We provide such a sufficient graphical criterion for \mathcal{D} -CI under an injectivity condition of a structural equation. Depending on the assumed causal graph, this can also come at the cost of requiring certain variables to be observed. As our main contribution, we propose a model-agnostic learning framework, called Counterfactually Invariant Prediction (CIP), using a kernel-based conditional dependence measure that also works for mixed categorical and continuous, multivariate variables. We evaluate CIP extensively in (semi-)synthetic settings and demonstrate its efficacy in enforcing counterfactual invariance even when the strict assumptions may be violated.

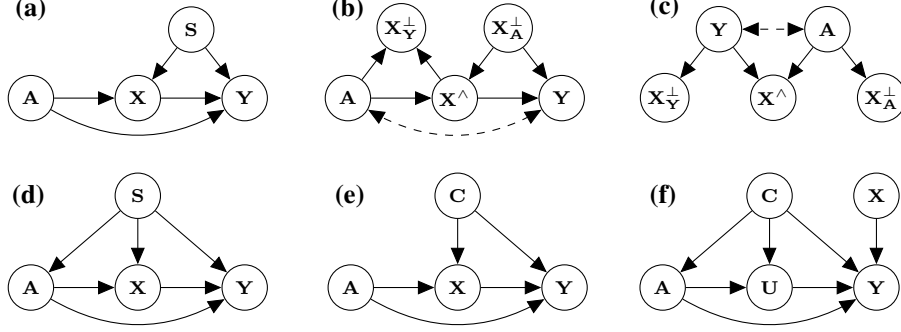


Figure 1: (a) Exemplary graph in which a predictor \hat{Y} with $\hat{Y} \perp\!\!\!\perp A \cup X \mid S$ is CI in A w.r.t. X . (b)-(c) Causal and anti-causal structure from Veitch et al. (2021) where X_A^\perp is not causally influenced by A , X_Y^\perp does not causally influence Y , and X^\wedge is both influenced by A and influences Y . (d) Assumed causal structure for the synthetic experiments, see Section 4.1 and Appendix F for details. (e) Assumed causal graph for the UCI Adult dataset (Section 4.3), where $A = \{\text{'gender', 'age'}\}$. (f) Causal structure for our semi-synthetic image experiments (Section 4.2), where $A = \{\text{Pos.X}\}$, $U = \{\text{Scale}\}$, $C = \{\text{Shape, Pos.Y}\}$, $X = \{\text{Color, Orientation}\}$, and $Y = \{\text{Outcome}\}$.

3 COUNTERFACTUALLY INVARIANT PREDICTION (CIP)

3.1 SUFFICIENT CRITERION FOR COUNTERFACTUAL INVARIANCE

We will now establish a sufficient graphical criterion to express CI as conditional independence in the observational distribution, rendering it estimable from data. First, we need some terminology.

Graph terminology. Consider a path π (a sequence of distinct adjacent nodes) in a DAG \mathcal{G} . A set of nodes S is said to *block* π , if π contains a triple of consecutive nodes A, B, C such that one of the following hold: (i) $A \rightarrow B \rightarrow C$ or $A \leftarrow B \leftarrow C$ or $A \leftarrow B \rightarrow C$ and $B \in S$; (ii) $A \rightarrow B \leftarrow C$ and neither B nor any descendant of B is in S . Further, we call π a *causal path* between sets of nodes A, B , when it is a directed path from a node in A to a node in B . A causal path π is a *proper causal path* if it only intersects A at the first node in π . Finally, we denote with \mathcal{G}_A the graph obtained by removing from \mathcal{G} all incoming arrows into nodes in A . We now define the notion of valid adjustment sets (Shpitser et al., 2010, Def. 5), which our graphical criterion for CI relies on.

Definition 3.1 (valid adjustment set). Let \mathcal{G} be a causal graph and let X, Y be disjoint (sets of) nodes in \mathcal{G} . A set of nodes S is a valid adjustment set for (X, Y) , if (i) no element in S is a descendant in \mathcal{G}_X of any node $W \notin X$ which lies on a proper causal path from X to Y , and (ii) S blocks all non-causal paths from X to Y in \mathcal{G} .

The sufficient graphical criterion that renders CI equivalent to a conditional independence then reads as follows. (The proof is deferred to Appendix A.)

Theorem 3.2. Let \mathcal{G} be a causal graph, A, W be two (not necessarily disjoint) sets of nodes in \mathcal{G} , such that $(A \cup W) \cap Y = \emptyset$, let S be a valid adjustment set for $(A \cup W, Y)$. Further, for $X := W \setminus A$ assume that $X = g(X, A, U_X)$ (which implies $\text{pa}(V) \in X \cup A$ for all $V \in X \cup A$) with g injective in U_X for all values of A and X . Then, in all SCMs compatible with \mathcal{G} , if a predictor \hat{Y} satisfies $\hat{Y} \perp\!\!\!\perp A \cup W \mid S$, then \hat{Y} is counterfactually invariant in A with respect to W .

Assumptions. First, we note that assuming the causal graph to be known is a standard assumption widely made in the causality literature, even though it is a strong one (Cartwright, 2007). Fawkes & Evans (2023, Prop. A.3) shows that CI cannot be decided from the observational distribution (and the causal graph) unless strong assumptions are made—they consider a case in which counterfactuals are identified. A strong assumption of Theorem 3.2, is on the arguments and injectivity of g , which ensures that any given observation (x, a) of X and A puts a point mass on a single u for U_X during abduction. While this is an untestable and admittedly strong, we highlight that this is not a limitation of our works specifically, but at least comparably strong assumptions are provably required for any method that claims to guarantee CI from observational data. Therefore, we complement our theoretical results by extensive experimental evaluation demonstrating CIPS’s efficacy even when some assumptions are violated. Finally, we do *not* assume the full SCM to be known.

3.2 EXAMPLE USE-CASES OF COUNTERFACTUALLY INVARIANT PREDICTION

Fig. 1(a) shows an exemplary graph in which the outcome Y is affected by (disjoint) sets A (in which we want to be CI) and X (inputs to f_Y). There may be confounders S between X and Y and we consider $W = X \cup A \cup S$. Here we aim to achieve $\hat{Y} \perp\!\!\!\perp A \cup X \mid S$ to obtain CI in A w.r.t. X . In our synthetic experiments, we also allow S to affect A , see Fig. 1(d). Let us further illustrate concrete potential applications of CI, which we later also study in our experiments.

Counterfactual fairness. Counterfactual fairness (Kusner et al., 2017) informally challenges a consequential decision: “*Would I have gotten the same outcome had my gender, race, or age been different with all else being equal?*”. Here $Y \subset V$ denotes the outcome and $A \subset V \setminus Y$ the *protected attributes* such as gender, race, or age—protected under anti-discrimination laws (Barocas & Selbst, 2016)—by $A \subset V \setminus Y$. Collecting all remaining observed covariates into $W := V \setminus Y$ counterfactual fairness reduces to counterfactual invariance. In experiments, we build a semi-synthetic DGP assuming the graph in Fig. 1(e) for the UCI adult dataset (Kohavi & Becker, 1996).

Robustness. CI serves as a strong notion of robustness in settings such as image classification: “*Would the truck have been classified correctly had it been winter in this situation instead of summer?*” For concrete demonstration, we use the dSprites dataset (Matthey et al., 2017) consisting of simple black and white images of different shapes (squares, ellipses, ...), sizes, orientations, and locations. We devise a DGP for this dataset with the graph depicted in Fig. 1(f).

Text classification. Veitch et al. (2021) motivate the importance of counterfactual invariance in text classification tasks. Specifically, they consider the causal and anti-causal structures depicted in Veitch et al. (2021, Fig. 1), which we replicate in Fig. 1(b,c). Both diagrams consist of protected attributes A , observed covariates X , and outcomes Y . To apply our sufficient criterion to their settings, we must assume that A and Y are unconfounded. We show in Appendix G.1 that CIP still performs on par with Veitch et al. (2021) even when this assumption is violated.

Theorem 3.2 provides a sufficient condition for CI (Definition 2.2) in terms of the conditional independence $\hat{Y} \perp\!\!\!\perp A \cup W \mid S$. We next develop an operator $\text{HSCIC}(\hat{Y}, A \cup W \mid S)$ that is (a) efficiently estimable from data, (b) differentiable, (c) a monotonic measure of conditional dependence, and (d) is zero if and only if $\hat{Y} \perp\!\!\!\perp A \cup W \mid S$. Hence, it is a practical objective to enforce CI.

3.3 HSCIC FOR CONDITIONAL INDEPENDENCE

Consider two sets of random variables Y and $A \cup W$, and denote with $(\Omega_Y, \mathcal{F}_Y)$ and $(\Omega_{A \cup W}, \mathcal{F}_{A \cup W})$ the respective measurable spaces. Suppose that we are given two RKHSs $\mathcal{H}_Y, \mathcal{H}_{A \cup W}$ over the support of Y and $A \cup W$ respectively. The tensor product space $\mathcal{H}_Y \otimes \mathcal{H}_{A \cup W}$ is defined as the space of functions of the form $(f \otimes g)(y, [a, w]) := f(y)g([a, w])$, for all $f \in \mathcal{H}_Y$ and $g \in \mathcal{H}_{A \cup W}$. The tensor product space yields a natural RKHS structure, with kernel k defined by $k(y \otimes [a, w], y' \otimes [a', w']) := k_Y(y, y')k_{A \cup W}([a, w], [a', w'])$. We refer the reader to Szabó & Sriperumbudur (2017) for more details on tensor product spaces.

Definition 3.3 (HSCIC). For (sets of) random variables $Y, A \cup W, S$, the HSCIC *between Y and $A \cup W$ given S* is defined as the real-valued random variable $\text{HSCIC}(Y, A \cup W \mid S) = H_{Y, A \cup W \mid S} \circ S$ where $H_{Y, A \cup W \mid S}$ is a real-valued deterministic function, defined as $H_{Y, A \cup W \mid S}(s) := \|\mu_{Y, A \cup W \mid S=s} - \mu_{Y \mid S=s} \otimes \mu_{A \cup W \mid S=s}\|$ with $\|\cdot\|$ the norm induced by the inner product of the tensor product space $\mathcal{H}_Y \otimes \mathcal{H}_{A \cup W}$.

Our Definition 3.3 is motivated by, but differs slightly from Park & Muandet (2020, Def. 5.3), which relies on the Bochner conditional expected value. While it is functionally equivalent (with the same implementation, see Eq. (2)), ours has the benefit of bypassing some technical assumptions required by Park & Muandet (2020) (see Appendices C and D for details). The HSCIC has the following important property, proved in Appendix B.

Theorem 3.4 (Theorem 5.4 by Park & Muandet (2020)). *If the kernel k of $\mathcal{H}_Y \otimes \mathcal{H}_{A \cup W}$ is characteristic⁶, $\text{HSCIC}(Y, A \cup W \mid S) = 0$ almost surely if and only if $Y \perp\!\!\!\perp A \cup W \mid S$.*

Because “most interesting” kernels such as the Gaussian and Laplacian kernels are characteristic, and the tensor product of translation-invariant characteristic kernels is characteristic again (Szabó & Sriperumbudur, 2017), this natural assumption is non-restrictive in practice. Combining Theorems 3.2 and 3.4, we can now use HSCIC to reliably achieve counterfactual invariance.

⁶The tensor product kernel k is characteristic if $\mathbb{P}_{Y, A \cup W} \mapsto \mathbb{E}_{y, [a, w]} [k(\cdot, y \otimes [a, w])]$ is injective.

Corollary 3.5. *Under the assumptions of Theorem 3.2, if $\text{HSCIC}(\hat{\mathbf{Y}}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}) = 0$ almost surely, then $\hat{\mathbf{Y}}$ is counterfactually invariant in \mathbf{A} with respect to \mathbf{W} .*

In addition, since HSCIC is defined in terms of the MMD (Definition 3.3 and Park & Muandet (2020, Def. 5.3)), it inherits the weak convergence property, i.e., if $\text{HSCIC}(\hat{\mathbf{Y}}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S})$ converges to zero, then the counterfactual distributions (for different intervention values \mathbf{a}) weakly converge to the same distribution. We refer to Simon-Gabriel & Schölkopf (2018); Simon-Gabriel et al. (2020) for a precise characterization. Hence, as HSCIC decreases, the predictor approaches counterfactual invariance and we need not drive HSCIC all the way to zero to obtain meaningful results.

3.4 LEARNING COUNTERFACTUALLY INVARIANT PREDICTORS (CIP)

Corollary 3.5 justifies our proposed objective, namely to minimize the following loss

$$\mathcal{L}_{\text{CIP}}(\hat{\mathbf{Y}}) = \mathcal{L}(\hat{\mathbf{Y}}) + \gamma \cdot \text{HSCIC}(\hat{\mathbf{Y}}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}), \quad [\text{CIP loss}] \quad (1)$$

where $\mathcal{L}(\hat{\mathbf{Y}})$ is a task-dependent loss function (e.g., cross-entropy for classification, or mean squared error for regression) and $\gamma \geq 0$ is a parameter that regulates the trade-off between predictive performance and counterfactual invariance.

The meaning of γ and how to choose it. The second term in Eq. (1) amounts to the additional objective of CI, which is typically at odds with predictive performance within the observational distribution \mathcal{L} . In practice, driving HSCIC to zero, i.e., viewing our task as a constrained optimization problem, typically deteriorates predictive performance too much to be useful for prediction—especially in small data settings.⁷ As the choice of γ amounts to choosing an operating point between predictive performance and CI, it cannot be selected in a data-driven fashion. As different settings call for different tradeoffs, we advocate for employing the following procedure: (i) Train an unconstrained predictor for a base predictive performance (e.g., 92% accuracy or 0.21 MSE). (ii) Fix a tolerance level α , indicating the maximally tolerable loss in predictive performance (e.g., at most 5% drop in accuracy or at most 10% increase in MSE). (iii) Perform a log-spaced binary search on γ (e.g., on $[1e-4, 1e4]$) to find the largest γ such that the predictive performance of the resulting predictor achieves predictive performance within the tolerance α —see Appendix F.6 for an illustration.

Estimating the HSCIC from samples. The key benefit of HSCIC as a conditional independence measure is that it does not require parametric assumptions on the underlying probability distributions, and it is applicable for any mixed, multi-dimensional data modalities, as long as we can define positive definite kernels on them. Given n samples $\{(\hat{\mathbf{y}}_i, \mathbf{a}_i, \mathbf{w}_i, \mathbf{s}_i)\}_{i=1}^n$, denote with $\hat{K}_{\hat{\mathbf{Y}}}$ the kernel matrix with entries $[\hat{K}_{\hat{\mathbf{Y}}}]_{i,j} := k_{\hat{\mathbf{Y}}}(\hat{\mathbf{y}}_i, \hat{\mathbf{y}}_j)$, and let $\hat{K}_{\mathbf{A} \cup \mathbf{W}}$ be the kernel matrix for $\mathbf{A} \cup \mathbf{W}$. We estimate the $H_{\hat{\mathbf{Y}}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}} \equiv H_{\hat{\mathbf{Y}}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}}(\cdot)$ as

$$\begin{aligned} \hat{H}_{\hat{\mathbf{Y}}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}}^2 &= \hat{w}_{\hat{\mathbf{Y}}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}}^T \left(\hat{K}_{\hat{\mathbf{Y}}} \odot \hat{K}_{\mathbf{A} \cup \mathbf{W}} \right) \hat{w}_{\hat{\mathbf{Y}}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}} \\ &\quad - 2 \left(\hat{w}_{\hat{\mathbf{Y}} \mid \mathbf{S}}^T \hat{K}_{\hat{\mathbf{Y}}} \hat{w}_{\hat{\mathbf{Y}}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}} \right) \left(\hat{w}_{\mathbf{A} \cup \mathbf{W} \mid \mathbf{S}}^T \hat{K}_{\mathbf{A} \cup \mathbf{W}} \hat{w}_{\hat{\mathbf{Y}}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}} \right) \\ &\quad + \left(\hat{w}_{\hat{\mathbf{Y}} \mid \mathbf{S}}^T \hat{K}_{\hat{\mathbf{Y}}} \hat{w}_{\hat{\mathbf{Y}} \mid \mathbf{S}} \right) \left(\hat{w}_{\mathbf{A} \cup \mathbf{W} \mid \mathbf{S}}^T \hat{K}_{\mathbf{A} \cup \mathbf{W}} \hat{w}_{\mathbf{A} \cup \mathbf{W} \mid \mathbf{S}} \right), \end{aligned} \quad (2)$$

where \odot is element-wise multiplication. The functions $\hat{w}_{\hat{\mathbf{Y}} \mid \mathbf{S}} \equiv \hat{w}_{\hat{\mathbf{Y}} \mid \mathbf{S}}(\cdot)$, $\hat{w}_{\mathbf{A} \cup \mathbf{W} \mid \mathbf{S}} \equiv \hat{w}_{\mathbf{A} \cup \mathbf{W} \mid \mathbf{S}}(\cdot)$, and $\hat{w}_{\hat{\mathbf{Y}}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}} \equiv \hat{w}_{\hat{\mathbf{Y}}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}}(\cdot)$ are found via kernel ridge regression. Caponnetto & Vito (2007) provide the convergence rates of the estimand $\hat{H}_{\hat{\mathbf{Y}}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}}^2$ under mild conditions. In practice, computing the HSCIC approximation by the formula in Eq. (2) can be computationally expensive. To speed it up, we can use random Fourier features to approximate $\hat{K}_{\hat{\mathbf{Y}}}$ and $\hat{K}_{\mathbf{A} \cup \mathbf{W}}$ (Rahimi & Recht, 2007; Avron et al., 2017). We emphasize that Eq. (2) allows us to consistently estimate the HSCIC from observational i.i.d. samples, without prior knowledge of the counterfactual distributions.

3.5 MEASURING COUNTERFACTUAL INVARIANCE.

Besides predictive performance, e.g., mean squared error (MSE) for regression or accuracy for classification, our key metric of interest is the level of counterfactual invariance achieved by the

⁷In particular, HSCIC does not regularize an ill-posed problem, i.e., it does not merely break ties between predictors with equal $\mathcal{L}(\hat{\mathbf{Y}})$. Hence it also need not decay to zero as the sample size increases.

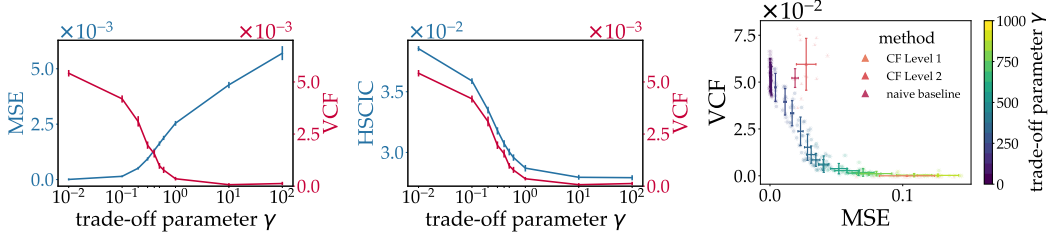


Figure 2: Results on synthetic data (see Appendix F.1 and Appendix F.2). **Left:** trade-off between MSE and counterfactual invariance (VCF). **Middle:** strong correspondence between HSCIC and VCF. **Right:** performance of CIP against baselines CF1 and CF2 and the naive baseline. As γ increases, CIP traces out a frontier characterizing the trade-off between MSE and CI. CF2 and the naive baseline are Pareto-dominated by CIP, i.e., we can pick γ to outperform CF2 in both MSE and VCF simultaneously. CF1 has zero VCF by design, but worse predictive performance than CIP at near zero VCF. Error bars are standard errors over 10 seeds.

predictor $\hat{\mathbf{Y}}$. First, we emphasize again that counterfactual distributions are generally not identified from the observational distribution (i.e., from available data) meaning that *CI is generally untestable in practice* from observational data. We can thus only evaluate CI in (semi-)synthetic settings where we have access to the full SCM and thus all counterfactual distributions.

A measure for CI must capture how the distribution of $\hat{\mathbf{Y}}_{\mathbf{a}'}$ changes for different values of \mathbf{a}' across all conditioning values \mathbf{w} (which may include an observed value $\mathbf{A} = \mathbf{a}$). We propose the **V**ariance of **C**ounter**F**actuals (VCF) as a metric of CI

$$\text{VCF}(\hat{\mathbf{Y}}) := \mathbb{E}_{\mathbf{w} \sim \mathbb{P}_{\mathbf{W}}} \left[\text{var}_{\mathbf{a}' \sim \mathbb{P}_{\mathbf{A}}} \left[\mathbb{E}_{\hat{\mathbf{Y}}_{\mathbf{a}'}}[\hat{\mathbf{Y}} | \mathbf{W} = \mathbf{w}] \right] \right]. \quad (3)$$

That is, we quantify how the average outcome varies with the interventional value \mathbf{a}' at conditioning value \mathbf{w} and average this variance over \mathbf{w} . For deterministic predictors (point estimators), which we use in all our experiments, the prediction is a fixed value for each input $\mathbb{E}_{\hat{\mathbf{Y}}_{\mathbf{a}'}}[\hat{\mathbf{Y}} | \mathbf{W} = \mathbf{w}] = \hat{\mathbf{y}}$ and we can drop the inner expectation of Eq. (3). In this case, the variance term in Eq. (3) is zero if and only if $\mathbb{P}_{\hat{\mathbf{Y}}_{\mathbf{a}'}}(\mathbf{y}) = \mathbb{P}_{\hat{\mathbf{Y}}_{\mathbf{a}'}}(\mathbf{y})$ almost surely. Since the variance is non-negative, the outer expectation is zero if and only if the variance term is zero almost surely, yielding the following result.

Corollary 3.6. *For point-estimators, $\hat{\mathbf{Y}}$ is counterfactually invariant in \mathbf{A} w.r.t. \mathbf{W} if and only if $\text{VCF}(\hat{\mathbf{Y}}) = 0$ almost surely.*

Estimating VCF in practice requires access to the DGP to generate counterfactuals. Given d i.i.d. examples $(\mathbf{w}_i)_{i=1}^d$ from a fixed observational dataset we sample k intervention values from the marginal $\mathbb{P}_{\mathbf{A}}$ and compute corresponding predictions. The inner expectation is simply the deterministic predictor output, and we compute the empirical expectation over the d observed \mathbf{w} values and empirical variances over the k sampled interventional values (for each \mathbf{w}). Since the required counterfactuals are by their very nature unavailable in practice, our analysis of VCF is limited to (semi-)synthetic settings. Notably, the proposed procedure for choosing γ does not require VCF. Our experiments corroborate the weak convergence property of HSCIC—small HSCIC implies small VCF. Hence, HSCIC may serve as a strong proxy for VCF and thus CI in practice.

4 EXPERIMENTS

Baselines. As many existing methods focus on cruder purely observational or interventional invariances (see Section 2.2), our choice of baselines for true counterfactual invariance is highly limited. First, we compare CIP to Veitch et al. (2021) in their two limited settings (Fig. 6(b-c)) in Appendix G.1, showing that our method performs on par with theirs. Next, we compare to two methods proposed by Kusner et al. (2017) in settings where they apply. CF1 (their ‘Level 1’) consists of only using non-descendants of \mathbf{A} as inputs to $f_{\hat{\mathbf{Y}}}$. CF2 (their ‘Level 2’) assumes an additive noise model and uses the residuals of descendants of \mathbf{A} after regression on \mathbf{A} together with non-descendants of \mathbf{A} as inputs to $f_{\hat{\mathbf{Y}}}$. We refer to these two baselines as CF1 and CF2 respectively. We also compare CIP to the ‘naive baseline’ which consists in training a predictor ignoring \mathbf{A} . In settings where \mathbf{A} is binary, we also compare to Chiappa (2019), devised for path-wise counterfactual fairness. Finally, we develop heuristics based on data augmentation as further possible baselines in Appendix G.2.

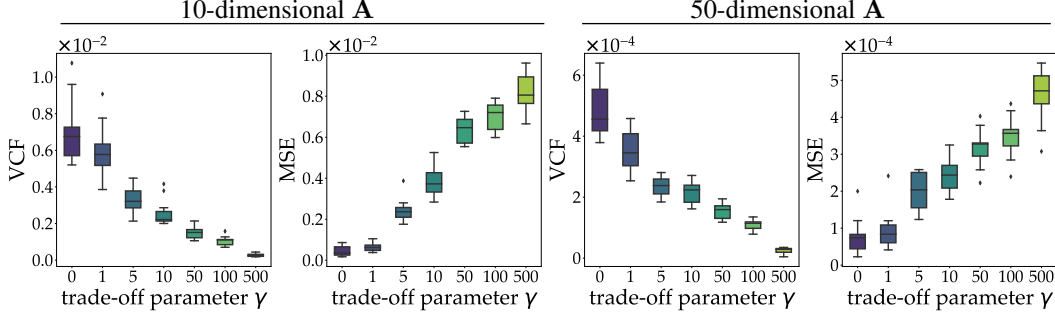


Figure 3: MSE and VCF for synthetic data (Appendix F.3) with 10- and 50-dimensional \mathbf{A} for different γ and 15 random seeds per box. CIP reliably achieves CI as γ increases.

4.1 SYNTHETIC EXPERIMENTS

First, we generate various synthetic datasets following the causal graph in Fig. 1(d). They contain (i) the prediction targets \mathbf{Y} , (ii) variable(s) we want to be CI in \mathbf{A} , (iii) covariates \mathbf{X} mediating effects from \mathbf{A} on \mathbf{Y} , and (iv) confounding variables \mathbf{S} . The goal is to learn a predictor $\hat{\mathbf{Y}}$ that is CI in \mathbf{A} w.r.t. $\mathbf{W} := \mathbf{A} \cup \mathbf{X} \cup \mathbf{S}$. The datasets cover different dimensions for the observed variables and their correlations and are described in detail in Appendix F.

Model choices and parameters. For all synthetic experiments, we train fully connected neural networks (MLPs) with MSE loss $\mathcal{L}_{\text{MSE}}(\hat{\mathbf{Y}})$ as the predictive loss \mathcal{L} in Eq. (1) for continuous outcomes \mathbf{Y} . We generate 10k samples from the observational distribution in each setting and use an 80 to 20 train-test split. All metrics reported are on the test set. We perform hyper-parameter tuning for MLP hyperparameters based on a random strategy (see Appendix F for details). The $\text{HSCIC}(\hat{\mathbf{Y}}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S})$ term is computed as in Eq. (2) using a Gaussian kernel with amplitude 1.0 and length scale 0.1. The regularization parameter λ for the ridge regression coefficients is set to $\lambda = 0.01$. We set $d = 1000$ and $k = 500$ in the estimation of VCF.

Model performance. We first study the effect of the HSCIC on accuracy and counterfactual invariance on the simulated dataset in Appendix F.1. Fig. 2 (left) depicts the expected trade-off between MSE and VCF for varying γ , whereas Fig. 2 (middle) highlights that HSCIC (estimable from observational data) is a strong proxy of counterfactual invariance measured by VCF (see discussion after Eq. (3)). Figure 2 (right) compares CIP to baselines for a simulated non-additive noise model in Appendix F.2. For a suitable choice of γ , CIP outperforms the baseline CF2 and the naive baseline in both MSE and VCF simultaneously. While CF1 achieves perfect CI by construction (VCF = 0), its MSE is higher than CIP at almost perfect CI (VCF near zero). To conclude, our method flexibly and reliably trades predictive performance for counterfactual invariance via a single parameter γ and Pareto-dominates existing methods. In Appendix F.2 we present extensive results on further simulated settings and compare CIP to other heuristic methods in Appendix G.2.

Effect of dimensionality of \mathbf{A} . A key advantage of CIP is that it can deal with multi-dimensional \mathbf{A} . We consider simulated datasets described in Appendix F.3, where we gradually increase the dimension of \mathbf{A} . The results in Fig. 3 for different trade-off parameters γ and different dimensions of \mathbf{A} demonstrate that CIP effectively enforces CI also for multi-dimensional \mathbf{A} .⁸ Further results are shown in Appendix F.3.

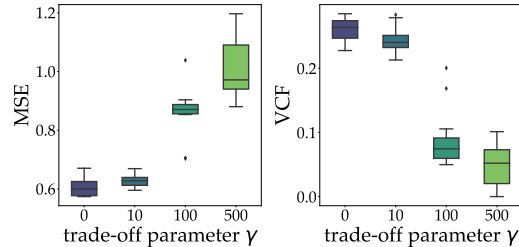


Figure 4: On the dSprites image dataset, CIP trades off MSE for VCF and achieves almost full CI as γ increases. Boxes are for 8 random seeds.

4.2 IMAGE EXPERIMENTS

We consider an image classification task on the dSprites dataset (Matthey et al., 2017), with a causal model as depicted in Fig. 1(f). The full structural equations are provided in Appendix F.4. This experiment is particularly challenging due to the mixed categorical and continuous variables in \mathbf{C} (shape, y-pos) and \mathbf{X} (color, orientation), with continuous \mathbf{A} (x-pos). We seek a

⁸In all boxplots, boxes represent the interquartile range, the horizontal line is the median, and whiskers show minimum and maximum values, excluding outliers (dots).

predictor \hat{Y} that is CI in the x-position w.r.t. all other observed variables. Following Theorem 3.2, we achieve $\hat{Y} \perp\!\!\!\perp \{x\text{-pos}, \text{scale}, \text{color}, \text{orientation}\} \mid \{\text{shape}, y\text{-pos}\}$ via the HSCIC operator. To accommodate the mixed input types, we first extract features from the images via a CNN and from other inputs via an MLP. We then use an MLP on all concatenated features for \hat{Y} . Fig. 4 shows that CIP gradually enforces CI as γ increases and illustrates the inevitable increase of MSE.

4.3 FAIRNESS WITH CONTINUOUS PROTECTED ATTRIBUTES

Finally, we apply CIP to the widely-used UCI Adult dataset (Kohavi & Becker, 1996) and compare it against a ‘naive baseline’ which simply ignores \mathbf{A} , CF1, CF2, and path-specific counterfactual fairness (PSCF) (Chiappa, 2019). We explicitly acknowledge the shortcomings of this dataset to reason about social justice (Ding et al., 2021). Instead, we chose it due to previous investigations into plausible causal structures based on domain knowledge (Zhang et al., 2017). The task is to predict whether an individual’s income is above a threshold based on demographic information, including protected attributes. We follow Nabi & Shpitser (2018); Chiappa (2019), where a causal structure is assumed for a subset of the variables as in Fig. 1(e) (see Appendix F.5 and Fig. 8 for details). We choose gender and age as the protected attributes \mathbf{A} , collect marital status, level of education, occupation, working hours per week, and work class into \mathbf{X} , and combine the remaining observed attributes in \mathbf{C} . Our aim is to learn a predictor \hat{Y} that is CI in \mathbf{A} w.r.t. $\mathbf{W} = \mathbf{C} \cup \mathbf{X}$. Achieving (causal) fairness for (mixed categorical and) continuous protected attributes is under active investigation (Mary et al., 2019; Chiappa & Pacchiano, 2021), but directly supported by CIP.

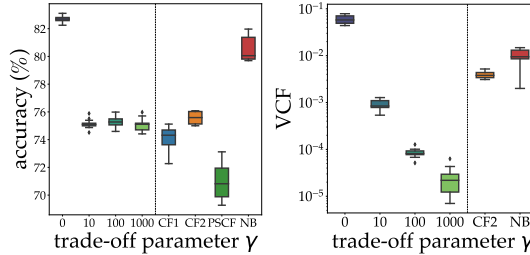


Figure 5: Accuracy and VCF on the Adult dataset. CIP achieves better VCF than CF2 and the naive baseline (NB), improved in accuracy compared to PSCF and is on par with CF1 in accuracy at $VCF \approx 0$.

We use an MLP with binary cross-entropy loss for \hat{Y} . Since this experiment is based on real data, the true counterfactual distribution cannot be known. Following Chiappa & Pacchiano (2021) we estimate a possible SCM by inferring the posterior distribution over the unobserved variables using variational autoencoders (Kingma & Welling, 2014). Figure 5 shows that CIP gradually achieves CI and even manages to keep a constant accuracy after an initial drop. It Pareto-dominates CF2 and PSCF and achieves comparable accuracy to CF1 when reaching $VCF \approx 0$. The naive baselines are more accurate than CIP for $\gamma \geq 5$ while CIP can achieve better VCF.

5 DISCUSSION AND FUTURE WORK

We developed CIP, a method to learn counterfactually invariant predictors \hat{Y} . First, we presented a sufficient graphical criterion to characterize counterfactual invariance and reduced it to conditional independence in the observational distribution under an injectivity assumption of a causal mechanism. We then built on kernel mean embeddings and the Hilbert-Schmidt Conditional Independence Criterion to devise an efficiently estimable, differentiable, model-agnostic objective to train CI predictors for mixed continuous/categorical, multi-dimensional variables. We demonstrated the efficacy of CIP in extensive empirical evaluations on various regression and classification tasks.

A key limitation of our work, shared by all studies in this domain, is the assumption that the causal graph is known. Guaranteeing CI necessarily requires strong untestable additional assumptions, but we demonstrated that CIP performs well empirically even when these are violated. Moreover, the computational cost of estimating HSCIC may limit the scalability of CIP to very high-dimensional settings even when using efficient random Fourier features. While the increased robustness of counterfactually invariant predictors are certainly desirable in many contexts, this presupposes the validity of our assumptions. Thus, an important direction for future work is to assess the sensitivity of CIP to misspecifications of the causal graph or insufficient knowledge of the required adjustment set. Lastly, we envision our graphical criterion and KME-based objective to be useful also for causal representation learning to isolate causally relevant, autonomous factors underlying the data.

REPRODUCIBILITY STATEMENT

The reproducibility of our work is ensured through several means. For the theoretical components, complete proofs of the stated theorems are provided in Appendix A and Appendix B. In terms of our experimental findings, Appendix F offers a detailed description of the hyperparameters and models. Moreover, to facilitate practical replication, the corresponding code is available in the supplementary material.

REFERENCES

- Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, 2019.
- Chen Avin, Ilya Shpitser, and Judea Pearl. Identifiability of path-specific effects. 2005.
- Haim Avron, Michael Kapralov, Cameron Musco, Christopher Musco, Ameya Velingker, and Amir Zandieh. Random fourier features for kernel ridge regression: Approximation bounds and statistical guarantees. In *International conference on machine learning*, volume 70, pp. 253–262, 2017.
- Elias Bareinboim, Juan D Correa, Duligur Ibeling, and Thomas Icard. On pearl’s hierarchy and the foundations of causal inference. In *Probabilistic and causal inference: the works of judea pearl*, pp. 507–556. 2022.
- Solon Barocas and Andrew D. Selbst. Big data’s disparate impact. *California Law Review*, 104, 2016.
- Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.
- Alain Berlinet and Christine Thomas-Agnan. *Reproducing kernel Hilbert spaces in probability and statistics*. Springer Science & Business Media, 2011.
- Benjamin Bloem-Reddy and Yee Whye Teh. Probabilistic symmetries and invariant neural networks. *The Journal of Machine Learning Research*, 21:90–1, 2020.
- Peter Bühlmann. Invariance, causality and robustness. *Statistical Science*, 35(3):404–426, 2020.
- Andrea Caponnetto and Ernesto De Vito. Optimal rates for the regularized least-squares algorithm. *Foundations of Computational Mathematics*, 7(3):331–368, 2007.
- Alycia N Carey and Xintao Wu. The causal fairness field guide: Perspectives from social and formal sciences. *Frontiers in Big Data*, 5:892837, 2022.
- Nancy Cartwright. *Hunting causes and using them: Approaches in philosophy and economics*. Cambridge University Press, 2007.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.
- Silvia Chiappa. Path-specific counterfactual fairness. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 7801–7808, 2019.
- Silvia Chiappa and Aldo Pacchiano. Fairness with continuous optimal transport. *arXiv preprint arXiv:2101.02084*, 2021.
- Erhan Çinlar and E. Çinlar. *Probability and stochastics*, volume 261. Springer, 2011.
- Taco Cohen and Max Welling. Group equivariant convolutional networks. In *International conference on machine learning*, pp. 2990–2999. PMLR, 2016.
- Nicolae Dinculeanu. *Vector integration and stochastic integration in Banach spaces*, volume 48. John Wiley & Sons, 2000.

- Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. Retiring adult: New datasets for fair machine learning. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- Cian Dorr. Against counterfactual miracles. *The Philosophical Review*, 125(2):241–286, 2016.
- Sanghamitra Dutta, Praveen Venkatesh, Piotr Mardziel, Anupam Datta, and Pulkrit Grover. Fairness under feature exemptions: Counterfactual and observational measures. *IEEE Transactions on Information Theory*, 67(10):6675–6710, 2021.
- Jake Fawkes and Robin J Evans. Results on counterfactual invariance. *arXiv preprint arXiv:2307.08519*, 2023.
- Kenji Fukumizu, Arthur Gretton, Xiaohai Sun, and Bernhard Schölkopf. Kernel measures of conditional dependence. In *Advances in neural information processing systems*, volume 20, 2007.
- Kenji Fukumizu, Le Song, and Arthur Gretton. Kernel bayes’ rule: Bayesian inference with positive definite kernels. *Journal of Machine Learning Research*, 14(1):3753–3783, 2013.
- Sainyam Galhotra, Karthikeyan Shanmugam, Prasanna Sattigeri, and Kush R Varshney. Causal feature selection for algorithmic fairness. In *Proceedings of the 2022 International Conference on Management of Data*, pp. 276–285, 2022.
- Arthur Gretton, Olivier Bousquet, Alexander J. Smola, and Bernhard Schölkopf. Measuring statistical dependence with hilbert-schmidt norms. In *Algorithmic Learning Theory*, volume 3734, pp. 63–77, 2005.
- Steffen Grünewälder, Guy Lever, Arthur Gretton, Luca Baldassarre, Sam Patterson, and Massimiliano Pontil. Conditional mean embeddings as regressors. In *International conference on machine learning*, 2012.
- Moritz Hardt, Eric Price, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In *Advances in Neural Information Processing Systems*, volume 29, 2016. URL <https://proceedings.neurips.cc/paper/2016/file/9d2682367c3935defcb1f9e247a97c0d-Paper.pdf>.
- Christina Heinze-Deml, Jonas Peters, and Nicolai Meinshausen. Invariant causal prediction for nonlinear models. *Journal of Causal Inference*, 6(2), 2018.
- Niki Kilbertus, Mateo Rojas Carulla, Giambattista Parascandolo, Moritz Hardt, Dominik Janzing, and Bernhard Schölkopf. Avoiding discrimination through causal reasoning. In *Advances in neural information processing systems*, volume 30, 2017.
- Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *2nd International Conference on Learning Representations, ICLR 2014*, 2014. URL <http://arxiv.org/abs/1312.6114>.
- Ronny Kohavi and Barry Becker. Uci adult data set. *UCI Machine Learning Repository*, 1996.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, volume 25, 2012. URL <https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf>.
- Matt J. Kusner, Joshua R. Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. In *Advances in Neural Information Processing Systems*, pp. 4066–4076, 2017.
- Joshua R Loftus, Chris Russell, Matt J Kusner, and Ricardo Silva. Causal reasoning for algorithmic fairness. *arXiv preprint arXiv:1805.05859*, 2018.
- Chaochao Lu, Yuhuai Wu, José Miguel Hernández-Lobato, and Bernhard Schölkopf. Invariant causal representation learning for out-of-distribution generalization. In *International Conference on Learning Representations*, 2021.

- Maggie Makar, Ben Packer, Dan Moldovan, Davis Blalock, Yoni Halpern, and Alexander D’Amour. Causally motivated shortcut removal using auxiliary labels. In *International Conference on Artificial Intelligence and Statistics*, pp. 739–766. PMLR, 2022.
- Jeremie Mary, Clément Calauzènes, and Nouredine El Karoui. Fairness-aware learning for continuous attributes and treatments. In *International Conference on Machine Learning*, volume 97, pp. 4382–4391. PMLR, 2019.
- Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner. dsprites: Disentanglement testing sprites dataset. <https://github.com/deepmind/dsprites-dataset/>, 2017.
- Shira Mitchell, Eric Potash, Solon Barocas, Alexander D’Amour, and Kristian Lum. Algorithmic fairness: Choices, assumptions, and definitions. *Annual Review of Statistics and Its Application*, 8: 141–163, 2021.
- S. Chandra Mouli and Bruno Ribeiro. Asymmetry learning for counterfactually-invariant classification in OOD tasks. In *Proc. of ICLR*, 2022.
- Krikamol Muandet, Kenji Fukumizu, Bharath Sriperumbudur, and Bernhard Schölkopf. Kernel mean embedding of distributions: A review and beyond. *Foundations and Trends in Machine Learning*, 10(1-2):1–141, 2017.
- Razieh Nabi and Ilya Shpitser. Fair inference on outcomes. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- Junhyung Park and Krikamol Muandet. A measure-theoretic approach to kernel conditional mean embeddings. In *Advances in Neural Information Processing Systems*, pp. 21247–21259, 2020.
- Judea Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, 2000.
- Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(5):947–1012, 2016.
- Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations and learning algorithms*. The MIT Press, 2017.
- Drago Plecko and Elias Bareinboim. Causal fairness analysis. *arXiv preprint arXiv:2207.11385*, 2022.
- Roman Pogodin, Namrata Deka, Yazhe Li, Danica J Sutherland, Victor Veitch, and Arthur Gretton. Efficient conditionally invariant representation learning. *arXiv preprint arXiv:2212.08645*, 2022.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems*, pp. 1177–1184, 2007.
- Mateo Rojas-Carulla, Bernhard Schölkopf, Richard Turner, and Jonas Peters. Invariant models for causal transfer learning. *The Journal of Machine Learning Research*, 19(1):1309–1342, 2018.
- Babak Salimi, Luke Rodriguez, Bill Howe, and Dan Suciu. Capuchin: Causal database repair for algorithmic fairness. *arXiv preprint arXiv:1902.08283*, 2019.
- Bernhard Schölkopf, Alexander J Smola, Francis Bach, et al. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press, 2002.
- Connor Shorten and Taghi M Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of big data*, 6(1):1–48, 2019.
- Ilya Shpitser and Judea Pearl. Complete identification methods for the causal hierarchy. *Journal of Machine Learning Research*, 9:1941–1979, 2008.
- Ilya Shpitser and Judea Pearl. Effects of treatment on the treated: Identification and generalization. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, pp. 514–521, 2009.

- Ilya Shpitser, Tyler J. 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*, pp. 527–536, 2010.
- Carl-Johann Simon-Gabriel and Bernhard Schölkopf. Kernel distribution embeddings: Universal kernels, characteristic kernels and kernel metrics on distributions. *Journal of Machine Learning Research*, 19(44):1–29, 2018.
- Carl-Johann Simon-Gabriel, Alessandro Barp, Bernhard Schölkopf, and Lester Mackey. Metrizing weak convergence with maximum mean discrepancies. *under review*, 2020.
- Alex Smola, Arthur Gretton, Le Song, and Bernhard Schölkopf. A hilbert space embedding for distributions. In *International Conference on Algorithmic Learning Theory*, pp. 13–31. Springer, 2007.
- Le Song, Jonathan Huang, Alexander J. Smola, and Kenji Fukumizu. Hilbert space embeddings of conditional distributions with applications to dynamical systems. In *International Conference on Machine Learning*, volume 382, pp. 961–968, 2009.
- Le Song, Kenji Fukumizu, and Arthur Gretton. Kernel embeddings of conditional distributions: A unified kernel framework for nonparametric inference in graphical models. *IEEE Signal Processing Magazine*, 30(4):98–111, 2013.
- Adarsh Subbaswamy, Bryant Chen, and Suchi Saria. A unifying causal framework for analyzing dataset shift-stable learning algorithms. *Journal of Causal Inference*, 10(1):64–89, 2022.
- Zoltán Szabó and Bharath K. Sriperumbudur. Characteristic and universal tensor product kernels. *Journal of Machine Learning Research*, 18:233:1–233:29, 2017.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Victor Veitch, Alexander D’Amour, Steve Yadlowsky, and Jacob Eisenstein. Counterfactual invariance to spurious correlations in text classification. In *Advances in Neural Information Processing Systems*, pp. 16196–16208, 2021.
- Julius von Kügelgen, Abdirisak Mohamed, and Sander Beckers. Backtracking counterfactuals. *arXiv preprint arXiv:2211.00472*, 2022.
- James Woodward. *Causation with a human face: Normative theory and descriptive psychology*. Oxford University Press, 2021.
- Qizhe Xie, Zihang Dai, Eduard Hovy, Thang Luong, and Quoc Le. Unsupervised data augmentation for consistency training. In *Advances in Neural Information Processing Systems*, volume 33, pp. 6256–6268, 2020. URL <https://proceedings.neurips.cc/paper/2020/file/44feb0096faa8326192570788b38c1d1-Paper.pdf>.
- Mingyang Yi, Ruoyu Wang, Jiacheng Sun, Zhenguo Li, and Zhi-Ming Ma. Breaking correlation shift via conditional invariant regularizer. In *The Eleventh International Conference on Learning Representations*, 2022.
- Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. In *Advances in Neural Information Processing Systems*, volume 30, 2017. URL <https://proceedings.neurips.cc/paper/2017/file/f22e4747dalaa27e363d86d40ff442fe-Paper.pdf>.
- Lu Zhang, Yongkai Wu, and Xintao Wu. Achieving non-discrimination in data release. In *International Conference on Knowledge Discovery and Data Mining*, pp. 1335–1344, 2017.

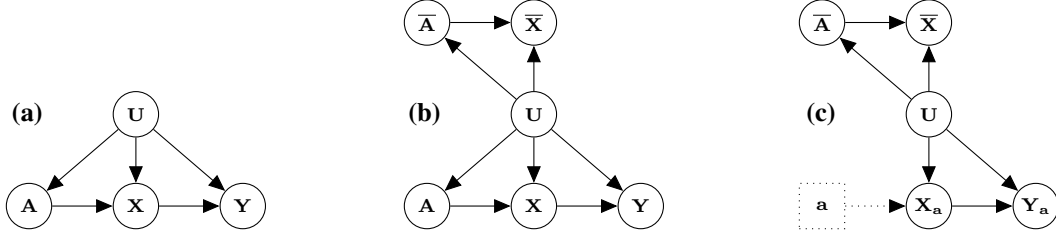


Figure 6: (a) A causal graph \mathcal{G} , which embeds information for the random variables of the model in the pre-interventional world. (b) The corresponding graph \mathcal{G}' for the set $\mathbf{W} = \{\mathbf{A}, \mathbf{X}\}$. The variables $\bar{\mathbf{A}}$ and $\bar{\mathbf{X}}$ are copies of \mathbf{A} and \mathbf{X} respectively. (c) The post-interventional graph \mathcal{G}'_a . By construction, any intervention of the form $\mathbf{A} \leftarrow \mathbf{a}$ does not affect the group $\bar{\mathbf{W}} = \{\bar{\mathbf{A}}, \bar{\mathbf{X}}\}$.

A PROOF OF THEOREM 3.2

A.1 OVERVIEW OF THE PROOF TECHNIQUES

We restate the main theorem for completeness.

Theorem 3.2. Let \mathcal{G} be a causal graph, \mathbf{A}, \mathbf{W} be two (not necessarily disjoint) sets of nodes in \mathcal{G} , such that $(\mathbf{A} \cup \mathbf{W}) \cap \mathbf{Y} = \emptyset$, let \mathbf{S} be a valid adjustment set for $(\mathbf{A} \cup \mathbf{W}, \mathbf{Y})$. Further, for $\mathbf{X} := \mathbf{W} \setminus \mathbf{A}$ assume that $\mathbf{X} = g(\mathbf{X}, \mathbf{A}, \mathbf{U}_{\mathbf{X}})$ (which implies $\text{pa}(V) \in \mathbf{X} \cup \mathbf{A}$ for all $V \in \mathbf{X} \cup \mathbf{A}$) with g injective in $\mathbf{U}_{\mathbf{X}}$ for all values of \mathbf{A} and \mathbf{X} . Then, in all SCMs compatible with \mathcal{G} , if a predictor $\hat{\mathbf{Y}}$ satisfies $\hat{\mathbf{Y}} \perp\!\!\!\perp \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}$, then $\hat{\mathbf{Y}}$ is counterfactually invariant in \mathbf{A} with respect to \mathbf{W} .

Our proof technique generalizes the work of Shpitser & Pearl (2009). To understand the proof technique, note that conditional counterfactual distributions of the form $\mathbb{P}_{\mathbf{Y}_{\mathbf{a}}^* | \mathbf{W}}(\mathbf{y} \mid \mathbf{w})$ involve quantities from two different worlds. The variables \mathbf{W} belong to the pre-interventional world, and the interventional variable $\mathbf{Y}_{\mathbf{a}}^*$ belongs to the world after performing the intervention $\mathbf{A} \leftarrow \mathbf{a}$. Hence, we study the identification of conditional counterfactual distributions using a diagram that embeds the causal relationships between the pre- and the post-interventional world. After defining this diagram, we prove that some conditional measures in this new model provide an estimate for $\mathbb{P}_{\mathbf{Y}_{\mathbf{a}}^* | \mathbf{W}}(\mathbf{y} \mid \mathbf{w})$. We then combine this result with the properties of \mathbf{Z} to prove the desired result.

A.2 IDENTIFIABILITY OF COUNTERFACTUAL DISTRIBUTIONS

In this section, we discuss a well-known criterion for the identifiability of conditional distributions, which we will then use to prove Theorem 3.2. To this end, we use the notions of a blocked path and valid adjustment set, which we restate for clarity.

Definition A.1. Consider a path π of causal graph \mathcal{G} . A set of nodes \mathbf{Z} blocks π , if π contains a triple of consecutive nodes connected in one of the following ways: $N_i \rightarrow Z \rightarrow N_j$, $N_i \leftarrow Z \rightarrow N_j$, with $N_i, N_j \notin \mathbf{Z}$, $Z \in \mathbf{Z}$, or $N_i \rightarrow M \leftarrow N_j$ and neither M nor any descendent of M is in \mathbf{Z} .

Using this definition, we define the concept of a valid adjustment set.

Definition 3.1 (valid adjustment set). Let \mathcal{G} be a causal graph and let \mathbf{X}, \mathbf{Y} be disjoint (sets of) nodes in \mathcal{G} . A set of nodes \mathbf{S} is a valid adjustment set for (\mathbf{X}, \mathbf{Y}) , if (i) no element in \mathbf{S} is a descendant in $\mathcal{G}_{\mathbf{X}}$ of any node $W \notin \mathbf{X}$ which lies on a proper causal path from \mathbf{X} to \mathbf{Y} , and (ii) \mathbf{S} blocks all non-causal paths from \mathbf{X} to \mathbf{Y} in \mathcal{G} .

Definition 3.1 is a useful graphical criterion for the identifiability of counterfactual distributions. In fact, following Corollary 1 by Shpitser et al. (2010), if \mathbf{S} satisfies the adjustment criterion relative to (\mathbf{A}, \mathbf{Y}) , then it holds

$$\mathbb{P}_{\mathbf{Y}_{\mathbf{a}}^*}(\mathbf{y}) = \int \mathbb{P}_{\mathbf{Y} | \mathbf{A}, \mathbf{S}}(\mathbf{y} \mid \mathbf{a}, \mathbf{s}) d\mathbb{P}_{\mathbf{S}}. \quad (4)$$

Furthermore, this identifiability criterion is *complete*. That is, consider any graph \mathcal{G} and a set of nodes \mathbf{S} that do not fulfill the valid adjustment criterion with respect to (\mathbf{A}, \mathbf{Y}) . Then, there exists a model inducing \mathcal{G} such that Eq. (4) does not hold (see Theorem 3 by Shpitser et al. (2010)).

A.3 d -SEPARATION AND CONDITIONAL INDEPENDENCE

In this section, we discuss a well-known criterion for conditional independence, which we will then use to prove Theorem 3.2. We use the notion of a blocked path, as in Definition A.2 and the concept of d -separation as follows.

Definition A.2 (d -Separation). Consider a causal graph \mathcal{G} . Two sets of nodes \mathbf{X} and \mathbf{Y} of \mathcal{G} are said to be d -separated by a third set \mathbf{S} if every path from any node of \mathbf{X} to any node of \mathbf{Y} is blocked by \mathbf{S} . We use the notation $\mathbf{X} \perp\!\!\!\perp_{\mathcal{G}} \mathbf{Y} \mid \mathbf{S}$ to indicate that \mathbf{X} and \mathbf{Y} are d -separated by \mathbf{S} in \mathcal{G} . We use Definition A.2 as a graphical criterion for conditional independence (Pearl, 2000).

Lemma A.3 (Markov Property). Consider a causal graph \mathcal{G} , and suppose that two sets of nodes \mathbf{X} and \mathbf{Y} of \mathcal{G} are d -separated by \mathbf{S} . Then, \mathbf{X} is independent of \mathbf{Y} given \mathbf{S} in any model induced by the graph \mathcal{G} .

The Markov Property is also referred to as d -separation property. We use the notation $\mathbf{X} \perp\!\!\!\perp_{\mathcal{G}} \mathbf{Y} \mid \mathbf{S}$ to indicate that \mathbf{X} and \mathbf{Y} are d -separated by \mathbf{S} in \mathcal{G} .

A.4 A GRAPHICAL CHARACTERIZATION OF COUNTERFACTUAL DISTRIBUTIONS

We study the relationships between the pre-interventional model corresponding to a causal diagram \mathcal{G} and the post-interventional model, inducing a diagram $\mathcal{G}_{\mathbf{a}}$ after an intervention $\mathbf{A} \leftarrow \mathbf{a}$. A natural way to study this relationship is to use the counterfactual graph (Shpitser & Pearl, 2008). However, the construction of the counterfactual graph is rather intricate. For our purposes it is sufficient to consider a simpler construction, generalizing the work by Shpitser & Pearl (2009).

Consider an SGM with causal graph \mathcal{G} , and fix a set of observed random variables of interest \mathbf{W} . Denote with $\text{de}(\mathbf{A})$ all descendants of \mathbf{A} in \mathcal{G} . Furthermore, for each node N of \mathcal{G} , denote with $\text{an}(N)$ the set of all its ancestral variables. We define the corresponding graph $\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}}$ in the following steps:

1. Define $\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}}$ to be the same graph as \mathcal{G} .
2. For each node $N \in \mathbf{A} \cup \mathbf{W}$, add a new duplicate node \bar{N} to $\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}}$.
3. For each node $N \in \mathbf{A} \cup \mathbf{W}$ and for each ancestral variable $P \in \text{an}(N) \setminus (\mathbf{A} \cup \mathbf{W})$ such that $P \in \text{de}((\mathbf{A} \cup \mathbf{W}))$, add a new duplicate node \bar{P} to $\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}}$.
4. For each duplicate node \bar{N} and for each parent $P \in \text{pa}(N)$, if a duplicate node \bar{P} was added in steps 2-3, then add an edge $\bar{P} \rightarrow \bar{N}$; otherwise add an edge $P \rightarrow \bar{N}$. The last part of this condition may be removed. In fact, assuming that it holds $\text{pa}(N) \subseteq \mathbf{A} \cup \mathbf{W}$, then a duplicate node \bar{P} was added in the previous steps, for any parent $P \in \text{pa}(N)$.
5. For each duplicate node \bar{N} , add an edge $U_N \rightarrow \bar{N}$.

An illustration of this graph is presented in Fig. 6. We denote with $\bar{\mathcal{H}}$ the set of duplicate nodes that were added to $\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}}$. We can naturally define structural equations for the new variables \bar{N} as

$$\bar{N} = f_N(\text{pa}(\bar{N}), U_N),$$

with f_N the structural equation for N in the original model, and $\text{pa}(\bar{N})$ the parents of \bar{N} in the newly define graph $\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}}$. Note that each random variable \bar{N} is a *copy* of the corresponding N , in the sense that $\bar{N} = N$ almost surely. Importantly, the following lemma holds.

Lemma A.4. Suppose that a set of nodes \mathbf{S} satisfies the adjustment criterion relative to $(\mathbf{A} \cup \mathbf{W}, \mathbf{Y})$ in \mathcal{G} . Then, \mathbf{S} satisfies the adjustment criterion relative to $(\mathbf{A} \cup \mathbf{W}, \mathbf{Y})$ in $\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}}$.

Proof. We prove the claim, by showing that all non-causal paths in $\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}}$ from $\mathbf{A} \cup \mathbf{W}$ to \mathbf{Y} are blocked by \mathbf{S} . Indeed, if \mathbf{S} satisfies the adjustment criterion relative to $(\mathbf{A} \cup \mathbf{W}, \mathbf{Y})$ in \mathcal{G} , then condition (i) of the adjustment criterion Definition 3.1 relative to $(\mathbf{A} \cup \mathbf{W}, \mathbf{Y})$ in $\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}}$ is satisfied. Let π be any such non-causal path in $\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}}$ from $\mathbf{A} \cup \mathbf{W}$ to \mathbf{Y} . If π does not cross any duplicate node, then it is blocked by \mathbf{S} . Otherwise, without loss of generality, we can decompose π in three paths, which we refer to as π_1 , π_2 , and π_3 . The path π_1 starts from a node in $\mathbf{A} \cup \mathbf{W}$ of \mathcal{G} , and it terminates in $\bar{\mathcal{H}}$. The path π_2 only contains nodes in a node in $\bar{\mathcal{H}}$, and the path π_3 starts from a node of $\bar{\mathcal{H}}$, and it terminates in \mathbf{Y} . The paths π_1 and π_3 necessarily contain paths of the form $\bar{N} \leftarrow P$

or $\bar{N} \leftarrow U_N \rightarrow N$, with $\bar{N} \in \bar{\mathcal{H}}$, P and N nodes of \mathcal{G} , and U_P a latent variable. By construction, no node $\bar{N} \in \bar{\mathcal{H}}$ belongs to the adjustment set \mathbf{S} . Hence, the path π contains a fork of three nodes, with the central node, or any descendants of the central node, are included in \mathbf{S} . Hence, the path π is blocked.

□

We further prove the following lemma.

Lemma A.5 (Following Theorem 4 by Shpitser et al. (2010)). *Define the sets $\mathbf{X} = \mathbf{W} \setminus \mathbf{A}$ and $\bar{\mathbf{X}} = \bar{\mathbf{W}} \setminus \bar{\mathbf{A}}$. Suppose that a set of nodes \mathbf{S} satisfies the adjustment criterion relative to $(\mathbf{A} \cup \mathbf{W}, \mathbf{Y})$ in \mathcal{G} . Then, it holds $\mathbf{Y}_{\mathbf{a}', \mathbf{x}'}^* \perp\!\!\!\perp \bar{\mathbf{A}}, \bar{\mathbf{X}} \mid \mathbf{S}$ for any intervention $\mathbf{A}, \mathbf{X} \leftarrow \mathbf{a}', \mathbf{x}'$.*

Proof. By Lemma A.4, if \mathbf{S} satisfies the adjustment criterion relative to $(\mathbf{A} \cup \mathbf{W}, \mathbf{Y})$ in \mathcal{G} , then it is also satisfies it in $\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}}$. Equivalently, \mathbf{S} satisfies the adjustment criterion relative to $(\mathbf{A} \cup \mathbf{X}, \mathbf{Y})$ in $\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}}$. Hence, by the sufficiency of the adjustment criterion (Theorem 4 by Shpitser et al. (2010)), it hold $\mathbf{Y} \perp\!\!\!\perp \mathbf{A}, \mathbf{X} \mid \mathbf{S}$ in the graph $(\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}})_{\mathbf{a}', \mathbf{x}'}$, which is obtained from $\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}}$ by performing an intervention $\mathbf{A}, \mathbf{X} \leftarrow \mathbf{a}', \mathbf{x}'$. By definition, the group of random variables $\bar{\mathbf{A}}$ and $\bar{\mathbf{X}}$ in $(\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}})_{\mathbf{a}', \mathbf{x}'}$ are copies of the pre-interventional variables \mathbf{A}, \mathbf{X} in $(\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}})_{\mathbf{a}', \mathbf{x}'}$. It follows that $\mathbf{Y} \perp\!\!\!\perp \bar{\mathbf{A}}, \bar{\mathbf{X}} \mid \mathbf{S}$ in the graph $(\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}})_{\mathbf{a}', \mathbf{x}'}$ or, equivalently, that $\mathbf{Y}_{\mathbf{a}', \mathbf{x}'}^* \perp\!\!\!\perp \bar{\mathbf{A}}, \bar{\mathbf{X}} \mid \mathbf{S}$, as claimed. □

A.5 PROOF OF THEOREM 3.2

We can identify conditional counterfactual distributions in \mathcal{G} , by identifying distributions on \mathcal{G}' . We can combine this observation with the notion of a valid adjustment set to derive a closed formula for the identification of the distributions of interest.

Before discussing the proof of Theorem 3.2, we prove an additional auxiliary result.

Lemma A.6. *Define the sets $\mathbf{X} = \mathbf{W} \setminus \mathbf{A}$ and $\bar{\mathbf{X}} = \bar{\mathbf{W}} \setminus \bar{\mathbf{A}}$. Then, for any intervention $\mathbf{A} \leftarrow \mathbf{a}'$ and observational values $\bar{\mathbf{a}}, \bar{\mathbf{x}}$, there exist an intervention $\bar{\mathbf{X}} \leftarrow \mathbf{x}'$ such that*

$$\bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}', \mathbf{x}'}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{y}) = \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}'}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{y})$$

Proof. Suppose that the following statement holds:

$$\text{there exist a point } \mathbf{x}' \text{ such that } \bar{\mathbb{P}}_{\mathbf{X}_{\mathbf{a}'} | \bar{\mathbf{X}}=\bar{\mathbf{x}}, \bar{\mathbf{A}}=\bar{\mathbf{a}}}(\mathbf{x}') = 1. \quad (5)$$

then our claim follows. In fact, for this point \mathbf{x}' it holds

$$\bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}', \mathbf{x}'}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}, \mathbf{X}_{\mathbf{a}'}=\mathbf{x}'}(\mathbf{y}) = \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}'}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}, \mathbf{X}_{\mathbf{a}'}=\mathbf{x}'}(\mathbf{y}) \quad (6)$$

Then, it holds

$$\begin{aligned} \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}', \mathbf{x}'}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{y}) &= \int \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}', \mathbf{x}'}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}, \mathbf{X}_{\mathbf{a}'}=\mathbf{x}'}(\mathbf{y}) d\bar{\mathbb{P}}_{\mathbf{X}_{\mathbf{a}'} | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{x}') \quad (\text{by Eq. (5)}) \\ &= \int \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}'}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}, \mathbf{X}_{\mathbf{a}'}=\mathbf{x}'}(\mathbf{y}) d\bar{\mathbb{P}}_{\mathbf{X}_{\mathbf{a}'} | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{x}') \quad (\text{by Eq. (6)}) \\ &= \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}'}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{y}). \quad (\text{by Eq. (5)}) \end{aligned}$$

Hence, the claim follows by showing that Eq. (5) holds.

To conclude the proof, we show that (5) holds. By our assumption, there exists a function g such that for each pair \mathbf{x}, \mathbf{a} there exist a unique point \mathbf{u} in the support of $U_{\mathbf{X}}$ such that $\mathbf{x} = g(\mathbf{x}, \mathbf{a}, \mathbf{u})$. By construction, it holds $\bar{\mathbf{X}} = g(\bar{\mathbf{X}}, \bar{\mathbf{A}}, U_{\mathbf{X}})$, from which it follows that for each pair $\bar{\mathbf{x}}, \bar{\mathbf{a}}$ there exist a unique point $\bar{\mathbf{u}}$ in the support of $U_{\mathbf{X}}$ such that $\bar{\mathbf{x}} = g(\bar{\mathbf{x}}, \bar{\mathbf{a}}, \bar{\mathbf{u}})$. Hence, for this point $\bar{\mathbf{u}}$ it holds $\bar{\mathbb{P}}_{U_{\bar{\mathbf{X}} | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}}(\bar{\mathbf{u}}) = 1$. It follows that it holds

$$\bar{\mathbb{P}}_{\mathbf{X}_{\mathbf{a}'} | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{x}') = 1 \quad \text{with} \quad \mathbf{x}' := g(\bar{\mathbf{x}}, \mathbf{a}', \bar{\mathbf{u}}).$$

Hence, (5) holds. □

We now prove our main result.

Proof of Theorem 3.2. Following the notation of Lemma A.5, define the sets $\mathbf{X} = \mathbf{W} \setminus \mathbf{A}$, $\bar{\mathbf{X}} = \bar{\mathbf{W}} \setminus \bar{\mathbf{A}}$, and let $\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}}$ be the augmented graph obtained by adding duplicate nodes. Note that, using this notation, the assumption that $\mathbf{Y} \perp\!\!\!\perp \mathbf{A}$, $\mathbf{W} \mid \mathbf{S}$ can be written as $\mathbf{Y} \perp\!\!\!\perp \mathbf{A}$, $\mathbf{X} \mid \mathbf{S}$. Denote with $\bar{\mathbb{P}}$ the induced measure on $\mathcal{G}'_{\mathbf{A} \cup \mathbf{W}}$. Suppose that it holds

$$\bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}', \mathbf{x}'}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{y}) = \int \mathbb{P}_{\mathbf{Y} | \mathbf{A}=\mathbf{a}', \mathbf{X}=\mathbf{x}', \mathbf{S}=\mathbf{s}}(\mathbf{y}) d\bar{\mathbb{P}}_{\mathbf{S} | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{s}) \quad (7)$$

for any intervention $\mathbf{A}, \mathbf{X} \leftarrow \mathbf{a}', \mathbf{x}'$, and for any possible value \mathbf{w} attained by \mathbf{W} . Assuming that Eq. (7) holds, we have that

$$\begin{aligned} \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}', \mathbf{x}'}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{y}) &= \int \mathbb{P}_{\mathbf{Y} | \mathbf{A}=\mathbf{a}', \mathbf{X}=\mathbf{x}', \mathbf{S}=\mathbf{s}}(\mathbf{y}) d\bar{\mathbb{P}}_{\mathbf{S} | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{z}) \quad (\text{assuming Eq. (7)}) \\ &= \int \mathbb{P}_{\mathbf{Y} | \mathbf{A}=\mathbf{a}, \mathbf{X}=\mathbf{x}, \mathbf{S}=\mathbf{s}}(\mathbf{y}) d\bar{\mathbb{P}}_{\mathbf{S} | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{s}) \quad (\mathbf{Y} \perp\!\!\!\perp \mathbf{A}, \mathbf{X} \mid \mathbf{S}) \\ &= \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}, \mathbf{x}''}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{y}), \quad (\text{assuming Eq. (7)}) \end{aligned} \quad (8)$$

for any intervention $\mathbf{X} \leftarrow \mathbf{x}''$. To conclude, define the set $\bar{\mathbf{T}} = \bar{\mathbf{A}} \setminus \bar{\mathbf{W}}$. It follows that

$$\begin{aligned} \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}', \mathbf{x}'}^* | \bar{\mathbf{W}}=\bar{\mathbf{w}}}(\mathbf{y}) &= \int \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}', \mathbf{x}'}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{y}) d\bar{\mathbb{P}}_{\bar{\mathbf{T}} | \bar{\mathbf{W}}=\bar{\mathbf{w}}}(\bar{\mathbf{t}}) \quad (\text{by conditioning}) \\ &= \int \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}, \mathbf{x}''}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{y}) d\bar{\mathbb{P}}_{\bar{\mathbf{T}} | \bar{\mathbf{W}}=\bar{\mathbf{w}}}(\bar{\mathbf{t}}) \quad (\text{by Eq. (8)}) \\ &= \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}, \mathbf{x}''}^* | \bar{\mathbf{W}}=\bar{\mathbf{w}}}(\mathbf{y}). \quad (\text{by unconditioning}) \end{aligned} \quad (9)$$

Since $\mathbf{X} \subseteq \mathbf{W}$, from the inequalities above it holds

$$\bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}'}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{y}) = \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}', \mathbf{x}'}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{y}) = \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}, \mathbf{x}''}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{y}) = \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{y}),$$

where the first and last inequality use Lemma A.6 for suitable values $\mathbf{x}', \mathbf{x}''$, and the equation in the middle follows from (9). The proof of Theorem 3.2 thus boils down to proving Eq. (7). To this end, we use the valid adjustment property of \mathbf{S} . Note that by Lemma A.5 it holds $\mathbf{Y}_{\mathbf{a}', \mathbf{x}'}^* \perp\!\!\!\perp \bar{\mathbf{A}}, \bar{\mathbf{X}} \mid \mathbf{S}$. Hence,

$$\begin{aligned} &\bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}', \mathbf{x}'}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{y}) \\ &= \int \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}', \mathbf{x}'}^* | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{a}}, \mathbf{S}=\mathbf{s}}(\mathbf{y}) d\bar{\mathbb{P}}_{\mathbf{S} | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{a}}}(\mathbf{s}) \quad (\text{by conditioning}) \\ &= \int \bar{\mathbb{P}}_{\mathbf{Y}_{\mathbf{a}', \mathbf{x}'}^* | \mathbf{S}=\mathbf{s}}(\mathbf{y}) d\bar{\mathbb{P}}_{\mathbf{S} | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{s}) \quad (\mathbf{Y}_{\mathbf{a}', \mathbf{x}'}^* \perp\!\!\!\perp \bar{\mathbf{A}}, \bar{\mathbf{X}} \mid \mathbf{S}) \\ &= \int \mathbb{P}_{\mathbf{Y} | \mathbf{A}=\mathbf{a}', \mathbf{X}=\mathbf{x}', \mathbf{S}=\mathbf{s}}(\mathbf{y}) d\bar{\mathbb{P}}_{\mathbf{S} | \bar{\mathbf{A}}=\bar{\mathbf{a}}, \bar{\mathbf{X}}=\bar{\mathbf{x}}}(\mathbf{s}), \quad (\text{by Lemma A.4}) \end{aligned}$$

and Eq. (7) follows. \square

B PROOF OF THEOREM 3.4

We prove that the HSCIC can be used to promote conditional independence, using a similar technique as Park & Muandet (2020). The following theorem holds.

Theorem 3.4 (Theorem 5.4 by Park & Muandet (2020)). *If the kernel k of $\mathcal{H}_{\mathbf{X}} \otimes \mathcal{H}_{\mathbf{A} \cup \mathbf{W}}$ is characteristic⁹, $\text{HSCIC}(\mathbf{Y}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}) = 0$ almost surely if and only if $\mathbf{Y} \perp\!\!\!\perp \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}$.*

Proof. By definition, we can write $\text{HSCIC}(\mathbf{Y}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}) = H_{\mathbf{Y}, \mathbf{A} \cup \mathbf{W} | \mathbf{S}} \circ \mathbf{S}$, where $H_{\mathbf{Y}, \mathbf{A} \cup \mathbf{W} | \mathbf{S}}$ is a real-valued deterministic function. Hence, the HSCIC is a real-valued random variable, defined over the same domain $\Omega_{\mathbf{S}}$ of the random variable \mathbf{X} .

⁹The tensor product kernel k is characteristic if $\mathbb{P}_{\mathbf{Y}, \mathbf{A} \cup \mathbf{W}} \mapsto \mathbb{E}_{\mathbf{y}, [\mathbf{a}, \mathbf{w}]} [k(\cdot, \mathbf{y} \otimes [\mathbf{a}, \mathbf{w}])]$ is injective.

We first prove that if $\text{HSCIC}(\mathbf{Y}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}) = 0$ almost surely, then it holds $\mathbf{Y} \perp\!\!\!\perp \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}$. To this end, consider an event $\Omega' \subseteq \Omega_{\mathbf{X}}$ that occurs almost surely, and such that it holds $(H_{\mathbf{Y}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{X}} \circ \mathbf{X})(\omega) = 0$ for all $\omega \in \Omega'$. Fix a sample $\omega \in \Omega'$, and consider the corresponding value $\mathbf{s}_\omega = \mathbf{S}(\omega)$, in the support of \mathbf{S} . It holds

$$\begin{aligned} \int k(\mathbf{y} \otimes [\mathbf{a}, \mathbf{w}], \cdot) d\mathbb{P}_{\mathbf{Y}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}=\mathbf{s}_\omega} &= \mu_{\mathbf{Y}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}=\mathbf{s}_\omega} && \text{(by definition)} \\ &= \mu_{\mathbf{Y} \mid \mathbf{S}=\mathbf{s}_\omega} \otimes \mu_{\mathbf{A} \cup \mathbf{W} \mid \mathbf{S}=\mathbf{s}_\omega} && \text{(since } \omega \in \Omega') \\ &= \int k_{\mathbf{Y}}(\mathbf{y}, \cdot) d\mathbb{P}_{\mathbf{Y} \mid \mathbf{S}=\mathbf{s}_\omega} \otimes \int k_{\mathbf{A} \cup \mathbf{W}}([\mathbf{a}, \mathbf{w}], \cdot) d\mathbb{P}_{\mathbf{A} \cup \mathbf{W} \mid \mathbf{S}=\mathbf{s}_\omega} && \text{(by definition)} \\ &= \int k_{\mathbf{Y}}(\mathbf{y}, \cdot) \otimes k_{\mathbf{A} \cup \mathbf{W}}([\mathbf{a}, \mathbf{w}], \cdot) d\mathbb{P}_{\mathbf{Y} \mid \mathbf{S}=\mathbf{s}_\omega} \mathbb{P}_{\mathbf{A} \cup \mathbf{W} \mid \mathbf{S}=\mathbf{s}_\omega}, && \text{(by Fubini's Theorem)} \end{aligned}$$

with $k_{\mathbf{Y}}$ and $k_{\mathbf{A} \cup \mathbf{W}}$ the kernels of $\mathcal{H}_{\mathbf{Y}}$ and $\mathcal{H}_{\mathbf{A} \cup \mathbf{W}}$ respectively. Since the kernel k of the tensor product space $\mathcal{H}_{\mathbf{Y}} \otimes \mathcal{H}_{\mathbf{A} \cup \mathbf{W}}$ is characteristic, then the kernels $k_{\mathbf{Y}}$ and $k_{\mathbf{A} \cup \mathbf{W}}$ are also characteristic. Hence, it holds $\mathbb{P}_{\mathbf{Y}, \mathbf{A} \mid \mathbf{S}=\mathbf{s}_\omega} = \mathbb{P}_{\mathbf{Y} \mid \mathbf{S}=\mathbf{s}_\omega} \mathbb{P}_{\mathbf{A} \mid \mathbf{S}=\mathbf{s}_\omega}$ for all $\omega \in \Omega'$. Since the event Ω' occurs almost surely, then $\mathbb{P}_{\mathbf{Y}, \mathbf{A} \mid \mathbf{S}=\mathbf{s}_\omega} = \mathbb{P}_{\mathbf{Y} \mid \mathbf{S}=\mathbf{s}_\omega} \mathbb{P}_{\mathbf{A} \mid \mathbf{S}=\mathbf{s}_\omega}$ almost surely, that is $\mathbf{Y} \perp\!\!\!\perp \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}$.

Assume now that $\mathbf{Y} \perp\!\!\!\perp \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}$. By definition there exists an event $\Omega'' \subseteq \Omega_{\mathbf{S}}$ such that $\mathbb{P}_{\mathbf{Y}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}=\mathbf{s}_\omega} = \mathbb{P}_{\mathbf{Y} \mid \mathbf{S}=\mathbf{s}_\omega} \mathbb{P}_{\mathbf{A} \cup \mathbf{W} \mid \mathbf{S}=\mathbf{s}_\omega}$ for all samples $\omega \in \Omega''$, with $\mathbf{s}_\omega = \mathbf{S}(\omega)$. It holds

$$\begin{aligned} \mu_{\mathbf{Y}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}=\mathbf{s}_\omega} &= \int k(\mathbf{y} \otimes [\mathbf{a}, \mathbf{w}], \cdot) d\mathbb{P}_{\mathbf{Y}, \mathbf{A} \cup \mathbf{W} \mid \mathbf{S}=\mathbf{s}_\omega} && \text{(by definition)} \\ &= \int k(\mathbf{y} \otimes [\mathbf{a}, \mathbf{w}], \cdot) d\mathbb{P}_{\mathbf{Y} \mid \mathbf{S}=\mathbf{s}_\omega} \mathbb{P}_{\mathbf{A} \cup \mathbf{W} \mid \mathbf{S}=\mathbf{s}_\omega} && \text{(since } \omega \in \Omega'') \\ &= \int k_{\mathbf{Y}}(\mathbf{y}, \cdot) k_{\mathbf{A} \cup \mathbf{W}}([\mathbf{a}, \mathbf{w}], \cdot) d\mathbb{P}_{\mathbf{Y} \mid \mathbf{S}=\mathbf{s}_\omega} \mathbb{P}_{\mathbf{A} \cup \mathbf{W} \mid \mathbf{S}=\mathbf{s}_\omega} && \text{(by definition of } k) \\ &= \int k_{\mathbf{Y}}(\mathbf{y}, \cdot) d\mathbb{P}_{\mathbf{Y} \mid \mathbf{S}=\mathbf{s}_\omega} \otimes \int k_{\mathbf{A} \cup \mathbf{W}}([\mathbf{a}, \mathbf{w}], \cdot) d\mathbb{P}_{\mathbf{A} \cup \mathbf{W} \mid \mathbf{S}=\mathbf{s}_\omega} && \text{(by Fubini's Theorem)} \\ &= \mu_{\mathbf{Y} \mid \mathbf{S}=\mathbf{s}_\omega} \otimes \mu_{\mathbf{A} \cup \mathbf{W} \mid \mathbf{S}=\mathbf{s}_\omega}. && \text{(by definition)} \end{aligned}$$

The claim follows. \square

C CONDITIONAL KERNEL MEAN EMBEDDINGS AND THE HSCIC

The notion of conditional kernel mean embeddings has already been studied in the literature. We show that, under stronger assumptions, our definition is equivalent to the definition by [Park & Muandet \(2020\)](#). In this section, without loss of generality we will assume that $\mathbf{W} = \emptyset$ and we will refer to the conditioning set as \mathbf{Z} .

C.1 CONDITIONAL KERNEL MEAN EMBEDDINGS AND CONDITIONAL INDEPENDENCE

We show that, under stronger assumptions, the HSCIC can be defined using the Bochner conditional expected value. The Bochner conditional expected value is defined as follows.

Definition C.1. Fix two random variables \mathbf{Y}, \mathbf{Z} taking value in a Banach space \mathcal{H} , and denote with $(\Omega, \mathcal{F}, \mathbb{P})$ their joint probability space. Then, the Bochner conditional expectation of \mathbf{Y} given \mathbf{Z} is any \mathcal{H} -valued random variable \mathbf{X} such that

$$\int_E \mathbf{Y} d\mathbb{P} = \int_E \mathbf{X} d\mathbb{P}$$

for all $E \in \sigma(\mathbf{Z}) \subseteq \mathcal{F}$, with $\sigma(\mathbf{Z})$ the σ -algebra generated by \mathbf{Z} . We denote with $\mathbb{E}[\mathbf{Y} \mid \mathbf{Z}]$ the Bochner expected value. Any random variable \mathbf{X} as above is a version of $\mathbb{E}[\mathbf{Y} \mid \mathbf{Z}]$.

The existence and almost sure uniqueness of the conditional expectation are shown in [Dinculeanu \(2000\)](#). Given a RKHS \mathcal{H} with kernel k over the support of \mathbf{Y} , [Park & Muandet \(2020\)](#) define the corresponding conditional kernel mean embedding as

$$\mu_{\mathbf{Y} \mid \mathbf{Z}} := \mathbb{E}[k(\cdot, \mathbf{y}) \mid \mathbf{Z}].$$

Note that, according to this definition, $\mu_{\mathbf{Y} \mid \mathbf{Z}}$ is an \mathcal{H} -valued random variable, not a single point of \mathcal{H} . [Park & Muandet \(2020\)](#) use this notion to define the HSCIC as follows.

Definition C.2 (The HSCIC according to [Park & Muandet \(2020\)](#)). Consider (sets of) random variables \mathbf{Y} , \mathbf{A} , \mathbf{Z} , and consider two RKHS $\mathcal{H}_{\mathbf{Y}}$, $\mathcal{H}_{\mathbf{A}}$ over the support of \mathbf{Y} and \mathbf{A} respectively. The HSCIC between \mathbf{Y} and \mathbf{A} given \mathbf{Z} is defined as the real-valued random variable

$$\omega \mapsto \|\mu_{\mathbf{Y}, \mathbf{A} | \mathbf{Z}}(\omega) - \mu_{\mathbf{Y} | \mathbf{Z}}(\omega) \otimes \mu_{\mathbf{A} | \mathbf{Z}}(\omega)\|,$$

for all samples ω in the domain $\Omega_{\mathbf{Z}}$ of \mathbf{Z} . Here, $\|\cdot\|$ the metric induced by the inner product of the tensor product space $\mathcal{H}_{\mathbf{Y}} \otimes \mathcal{H}_{\mathbf{Z}}$.

We show that, under more restrictive assumptions, Definition C.2 can be used to promote conditional independence. To this end, we use the notion of a regular version.

Definition C.3 (Regular Version, following Definition 2.4 by [Çınlar & Şınlar \(2011\)](#)). Consider two random variables \mathbf{Y} , \mathbf{Z} , and consider the induced measurable spaces $(\Omega_{\mathbf{Y}}, \mathcal{F}_{\mathbf{Y}})$ and $(\Omega_{\mathbf{Z}}, \mathcal{F}_{\mathbf{Z}})$. A regular version Q for $\mathbb{P}_{\mathbf{Y} | \mathbf{Z}}$ is a mapping $Q: \Omega_{\mathbf{Z}} \times \mathcal{F}_{\mathbf{Y}} \rightarrow [0, +\infty]: (\omega, \mathbf{y}) \mapsto Q_{\omega}(\mathbf{y})$ such that: (i) the map $\omega \mapsto Q_{\omega}(\mathbf{x})$ is $\mathcal{F}_{\mathbf{A}}$ -measurable for all \mathbf{y} ; (ii) the map $\mathbf{y} \mapsto Q_{\omega}(\mathbf{y})$ is a measure on $(\Omega_{\mathbf{Y}}, \mathcal{F}_{\mathbf{Y}})$ for all ω ; (iii) the function $Q_{\omega}(\mathbf{y})$ is a version for $\mathbb{E} [\mathbb{1}_{\{\mathbf{Y}=\mathbf{y}\}} | \mathbf{Z}]$.

The following theorem shows that the random variable as in Definition C.2 can be used to promote conditional independence.

Theorem C.4 (Theorem 5.4 by [Park & Muandet \(2020\)](#)). *With the notation introduced above, suppose that the kernel k of the tensor product space $\mathcal{H}_{\mathbf{X}} \otimes \mathcal{H}_{\mathbf{A}}$ is characteristic. Furthermore, suppose that $\mathbb{P}_{\mathbf{Y}, \mathbf{A} | \mathbf{X}}$ admits a regular version. Then, $\|\mu_{\mathbf{Y}, \mathbf{A} | \mathbf{Z}}(\omega) - \mu_{\mathbf{Y} | \mathbf{Z}}(\omega) \otimes \mu_{\mathbf{A} | \mathbf{Z}}(\omega)\| = 0$ almost surely if and only if $\mathbf{Y} \perp\!\!\!\perp \mathbf{A} | \mathbf{Z}$.*

Note that the assumption of the existence of a regular version is essential in Theorem C.4. In this work, HSCIC is not used for conditional independence testing but as a conditional independence measure.

C.2 EQUIVALENCE WITH OUR APPROACH

The following theorem shows that under the existence of a regular version, conditional kernel mean embeddings can be defined using the Bochner conditional expected value. To this end, we use the following theorem.

Theorem C.5 (Following Proposition 2.5 by [Çınlar & Şınlar \(2011\)](#)). *Following the notation introduced in Definition C.3, suppose that $\mathbb{P}_{\mathbf{Y} | \mathbf{Z}}(\cdot | \mathbf{Z})$ admits a regular version $Q_{\omega}(\mathbf{y})$. Consider a kernel k over the support of \mathbf{Y} . Then, the mapping*

$$\omega \mapsto \int k(\cdot, \mathbf{y}) dQ_{\omega}(\mathbf{y})$$

is a version of $\mathbb{E} [k(\cdot, \mathbf{y}) | \mathbf{Z}]$.

As a consequence of Theorem C.5, we prove the following result.

Lemma C.6. *Fix two random variables \mathbf{Y} , \mathbf{Z} . Suppose that $\mathbb{P}_{\mathbf{Y} | \mathbf{Z}}$ admits a regular version. Denote with $\Omega_{\mathbf{Z}}$ the domain of \mathbf{Z} . Then, there exists a subset $\Omega \subseteq \Omega_{\mathbf{Z}}$ that occurs almost surely, such that $\mu_{\mathbf{Y} | \mathbf{Z}}(\omega) = \mu_{\mathbf{Y} | \mathbf{Z}=\mathbf{Z}(\omega)}$ for all $\omega \in \Omega$. Here, $\mu_{\mathbf{Y} | \mathbf{Z}=\mathbf{Z}(\omega)}$ is the embedding of conditional measures as in Section 2.*

Proof. Let $Q_{\omega}(\mathbf{y})$ be a regular version of $\mathbb{P}_{\mathbf{Y} | \mathbf{Z}}$. Without loss of generality we may assume that it holds $\mathbb{P}_{\mathbf{Y} | \mathbf{Z}}(\mathbf{y} | \{\mathbf{Z} = \mathbf{Z}(\omega)\}) = Q_{\omega}(\mathbf{y})$. By Theorem C.5 there exists an event $\Omega \subseteq \Omega_{\mathbf{Z}}$ that occurs almost surely such that

$$\mu_{\mathbf{Y} | \mathbf{Z}}(\omega) = \mathbb{E}[k(\mathbf{y}, \cdot) | \mathbf{Z}](\omega) = \int k(\mathbf{y}, \cdot) dQ_{\omega}(\mathbf{y}), \quad (10)$$

for all $\omega \in \Omega$. Then, for all $\omega \in \Omega$ it holds

$$\begin{aligned} \mu_{\mathbf{Y} | \mathbf{Z}}(\omega) &= \int k(\mathbf{x}, \cdot) dQ_{\omega}(\mathbf{x}) && \text{(it follows from Eq. (10))} \\ &= \int k(\mathbf{x}, \cdot) d\mathbb{P}_{\mathbf{X} | \mathbf{A}}(\mathbf{x} | \{\mathbf{A} = \mathbf{A}(\omega)\}) && (Q_{\omega}(\mathbf{y}) = \mathbb{P}_{\mathbf{Y} | \mathbf{Z}}(\mathbf{y} | \{\mathbf{Z} = \mathbf{Z}(\omega)\})) \\ &= \mu_{\mathbf{X} | \{\mathbf{A}=\mathbf{A}(\omega)\}}, && \text{(by definition as in Section 2)} \end{aligned}$$

as claimed. \square

As a consequence of Lemma C.6, we can prove that the definition of the HSCIC by Park & Muandet (2020) is equivalent to ours. The following corollary holds.

Corollary C.7. *Consider (sets of) random variables \mathbf{Y} , \mathbf{A} , \mathbf{Z} , and consider two RKHS $\mathcal{H}_{\mathbf{Y}}$, $\mathcal{H}_{\mathbf{A}}$ over the support of \mathbf{Y} and \mathbf{A} respectively. Suppose that $\mathbb{P}_{\mathbf{Y},\mathbf{A}|\mathbf{Z}}(\cdot | \mathbf{Z})$ admits a regular version. Then, there exists a set $\Omega \subseteq \Omega_{\mathbf{A}}$ that occurs almost surely, such that*

$$\|\mu_{\mathbf{X},\mathbf{A}|\mathbf{Z}}(\omega) - \mu_{\mathbf{X}|\mathbf{Z}}(\omega) \otimes \mu_{\mathbf{A}|\mathbf{Z}}(\omega)\| = (H_{\mathbf{Y},\mathbf{A}|\mathbf{Z}} \circ \mathbf{Z})(\omega).$$

Here, $H_{\mathbf{Y},\mathbf{A}|\mathbf{Z}}$ is a real-valued deterministic function, defined as

$$H_{\mathbf{Y},\mathbf{A}|\mathbf{Z}}(\mathbf{z}) := \|\mu_{\mathbf{Y},\mathbf{A}|\mathbf{Z}=\mathbf{z}} - \mu_{\mathbf{Y}|\mathbf{Z}=\mathbf{z}} \otimes \mu_{\mathbf{A}|\mathbf{Z}=\mathbf{z}}\|,$$

and $\|\cdot\|$ is the metric induced by the inner product of the tensor product space $\mathcal{H}_{\mathbf{X}} \otimes \mathcal{H}_{\mathbf{A}}$.

We remark that the assumption of the existence of a regular version is essential in Corollary C.7.

D THE CROSS-COVARIANCE OPERATOR

In this section, we show that under additional assumptions, our definition of conditional KMEs is equivalent to the definition based on the cross-covariance operator, under more restrictive assumptions. The definition of KMEs based on the cross-covariance operator requires the use of the following well-known result.

Lemma D.1. *Fix two RKHS $\mathcal{H}_{\mathbf{X}}$ and $\mathcal{H}_{\mathbf{Z}}$, and let $\{\varphi_i\}_{i=1}^{\infty}$ and $\{\psi_j\}_{j=1}^{\infty}$ be orthonormal bases of $\mathcal{H}_{\mathbf{X}}$ and $\mathcal{H}_{\mathbf{Z}}$ respectively. Denote with $\text{HS}(\mathcal{H}_{\mathbf{X}}, \mathcal{H}_{\mathbf{Z}})$ the set of Hilbert-Schmidt operators between $\mathcal{H}_{\mathbf{X}}$ and $\mathcal{H}_{\mathbf{Z}}$. There is an isometric isomorphism between the tensor product space $\mathcal{H}_{\mathbf{X}} \otimes \mathcal{H}_{\mathbf{Z}}$ and $\text{HS}(\mathcal{H}_{\mathbf{X}}, \mathcal{H}_{\mathbf{Z}})$, given by the map*

$$T: \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} c_{i,j} \varphi_i \otimes \psi_j \mapsto \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} c_{i,j} \langle \cdot, \varphi_i \rangle_{\mathcal{H}_{\mathbf{X}}} \psi_j.$$

For proof of this result see i.e., Park & Muandet (2020). This lemma allows us to define the cross-covariance operator between two random variables, using the operator T .

Definition D.2 (Cross-Covariance Operator). Consider two random variables \mathbf{X} , \mathbf{Z} . Consider corresponding mean embeddings $\mu_{\mathbf{X},\mathbf{Z}}$, $\mu_{\mathbf{X}}$ and $\mu_{\mathbf{Z}}$, as defined in Section 3. The cross-covariance operator is defined as $\Sigma_{\mathbf{X},\mathbf{Z}} := T(\mu_{\mathbf{X},\mathbf{Z}} - \mu_{\mathbf{X}} \otimes \mu_{\mathbf{Z}})$. Here, T is the isometric isomorphism as in Lemma D.1.

It is well-known that the cross-covariance operator can be decomposed into the covariance of the marginals and the correlation. That is, there exists a unique bounded operator $\Lambda_{\mathbf{Y},\mathbf{Z}}$ such that

$$\Sigma_{\mathbf{Y},\mathbf{Z}} = \Sigma_{\mathbf{Y},\mathbf{Y}}^{1/2} \circ \Lambda_{\mathbf{Y},\mathbf{Z}} \circ \Sigma_{\mathbf{Z},\mathbf{Z}}^{1/2}$$

Using this notation, we define the *normalized conditional cross-covariance operator*. Given three random variables \mathbf{Y} , \mathbf{A} , \mathbf{Z} and corresponding kernel mean embeddings, this operator is defined as

$$\Lambda_{\mathbf{Y},\mathbf{A}|\mathbf{Z}} := \Lambda_{\mathbf{Y},\mathbf{A}} - \Lambda_{\mathbf{Y},\mathbf{Z}} \circ \Lambda_{\mathbf{Z},\mathbf{A}}. \quad (11)$$

This operator was introduced by Fukumizu et al. (2007). The normalized conditional cross-covariance can be used to promote statistical independence, as shown in the following theorem.

Theorem D.3 (Theorem 3 by Fukumizu et al. (2007)). *Following the notation introduced above, define the random variable $\tilde{\mathbf{A}} := (\mathbf{A}, \mathbf{Z})$. Let $\mathbb{P}_{\mathbf{Z}}$ be the distribution of the random variable \mathbf{Z} , and denote with $L^2(\mathbb{P}_{\mathbf{Z}})$ the space of the square integrable functions with probability $\mathbb{P}_{\mathbf{Z}}$. Suppose that the tensor product kernel $k_{\mathbf{Y}} \otimes k_{\mathbf{A}} \otimes k_{\mathbf{Z}}$ is characteristic. Furthermore, suppose that $\mathcal{H}_{\mathbf{Z}} + \mathbb{R}$ is dense in $L^2(\mathbb{P}_{\mathbf{Z}})$. Then, it holds*

$$\Lambda_{\mathbf{Y},\tilde{\mathbf{A}}|\mathbf{Z}} = 0 \quad \text{if and only if} \quad \mathbf{Y} \perp\!\!\!\perp \mathbf{A} | \mathbf{Z}.$$

Here, $\Lambda_{\mathbf{Y},\tilde{\mathbf{A}}|\mathbf{Z}}$ is an operator defined as in Eq. (11).

By Theorem D.3, the operator $\Lambda_{\mathbf{Y},\tilde{\mathbf{A}}|\mathbf{Z}}$ can also be used to promote conditional independence. However, CIP is more straightforward since it requires less assumptions. In fact, Theorem D.3 requires to embed the variable \mathbf{Z} in an RKHS. In contrast, CIP only requires the embedding of the variables \mathbf{Y} and \mathbf{A} .

E RANDOM FOURIER FEATURES

Random Fourier features is an approach to scaling up kernel methods for shift-invariant kernels (Rahimi & Recht, 2007). Recall that a shift-invariant kernel is a kernel of the form $k(\mathbf{z}, \mathbf{z}') = h_k(\mathbf{z} - \mathbf{z}')$, with h_k a positive definite function.

Fourier features are defined via the following well-known theorem.

Theorem E.1 (Bochner’s Theorem). *For every shift-invariant kernel of the form $k(\mathbf{z}, \mathbf{z}') = h_k(\mathbf{z} - \mathbf{z}')$ with $h_k(\mathbf{0}) = 1$, there exists a probability density function $\mathbb{P}_k(\boldsymbol{\eta})$ such that*

$$k(\mathbf{z}, \mathbf{z}') = \int e^{-2\pi i \boldsymbol{\eta}^T (\mathbf{z} - \mathbf{z}')} d\mathbb{P}_k.$$

Since both the kernel k and the probability distribution \mathbb{P}_k are real-valued functions, the integrand in Theorem E.1 can be replaced by the function $\cos \boldsymbol{\eta}^T (\mathbf{z} - \mathbf{z}')$, and we obtain the following formula

$$k(\mathbf{z}, \mathbf{z}') = \int \cos \boldsymbol{\eta}^T (\mathbf{z} - \mathbf{z}') d\mathbb{P}_k = \mathbb{E} [\cos \boldsymbol{\eta}^T (\mathbf{z} - \mathbf{z}')] , \quad (12)$$

where the expected value is taken with respect to the distribution $\mathbb{P}_k(\boldsymbol{\eta})$. This equation allows to approximate the kernel $k(\mathbf{z}, \mathbf{z}')$, via the empirical mean of points $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_l$ sampled independently according to \mathbb{P}_k . In fact, it is possible to prove exponentially fast convergence of an empirical estimate for $\mathbb{E} [\cos \boldsymbol{\eta}^T (\mathbf{z} - \mathbf{z}')]$, as shown in the following theorem.

Theorem E.2 (Uniform Convergence of Fourier Features, Claim 1 by Rahimi & Recht (2007)). *Following the notation introduced above, fix any compact subset Ω in the domain of k , and consider points $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_l$ sampled independent according to the distribution \mathbb{P}_k . Define the function*

$$\hat{k}(\mathbf{z}, \mathbf{z}') := \frac{1}{l} \sum_{j=1}^l \cos \boldsymbol{\eta}_j^T (\mathbf{z} - \mathbf{z}'),$$

for all $(\mathbf{z}, \mathbf{z}') \in \Omega$. Then, it holds

$$\mathbb{P} \left(\sup_{\mathbf{z}, \mathbf{z}'} |\hat{k}(\mathbf{z}, \mathbf{z}') - k(\mathbf{z}, \mathbf{z}')| \geq \varepsilon \right) \leq 2^8 \sigma_k \frac{\text{diam}(\Omega)}{\varepsilon} \exp \left\{ -\frac{\varepsilon^2 l}{4(d+1)} \right\}.$$

Here σ_k^2 is the second moment of the Fourier transform of the kernel k , and d is the dimension of the arrays \mathbf{z} and \mathbf{z}' .

By Theorem E.2, the estimated kernel \hat{k} is a good approximation of the true kernel k on the set Ω .

Similarly, we can approximate the Kernel matrix using Random Fourier features. Following the notation introduced above, define the function

$$\zeta_{k,l}(\mathbf{z}) := \frac{1}{\sqrt{l}} [\cos \boldsymbol{\eta}_1^T \mathbf{z}, \dots, \cos \boldsymbol{\eta}_l^T \mathbf{z}] \quad (13)$$

with $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_l$ sampled independent according to the distribution \mathbb{P}_k .

We can approximate the Kernel matrix using the functions defined as in Eq. (13). Consider n samples $\mathbf{z}_1, \dots, \mathbf{z}_n$, and denote with Z the $n \times l$ matrix whose i -th row is given by $\zeta_{k,l}(\mathbf{z}_i)$. Similarly, denote with Z^* the $l \times n$ matrix whose i -th column is given by $\zeta_{k,l}^*(\mathbf{z}_i)$. Then, we can approximate the kernel matrix as $\hat{K}_Z \approx ZZ^*$.

We can also use this approximation to compute the kernel ridge regression parameters as in Section 3 using the formula $\hat{w}_{\mathbf{Y}|\mathbf{Z}}(\cdot) \approx (ZZ^* - n\lambda I)^{-1} [k_{\mathbf{Z}}(\cdot, \mathbf{z}_1), \dots, k_{\mathbf{Z}}(\cdot, \mathbf{z}_n)]^T$. Avron et al. (2017) argue that the approximate kernel ridge regression, as defined above, is an accurate estimate of the true distribution. Their argument is based on proving that the matrix $ZZ^* - n\lambda I$ is a good approximation of $\hat{K}_Z - n\lambda I$. The notion of good approximation is clarified by the following definition.

Definition E.3. Fix two Hermitian matrices A and B of the same size. We say that a matrix A is a γ -spectral approximation of another matrix B , if it holds $(1 - \gamma)B \preceq A \preceq (1 + \gamma)B$. Here, the \preceq symbol means that $A - (1 - \gamma)B$ is positive semi-definite, and that $(1 + \gamma)B - A$ is positive semi-definite.

Avron et al. (2017) prove that $ZZ^* - n\lambda I$ is a γ -approximation of $\hat{K}_{\mathbf{Z}} - n\varepsilon I$, if the number of samples $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_l$ is sufficiently large.

Theorem E.4 (Theorem 7 by Avron et al. (2017)). *Fix a constant $\gamma \leq 1/2$. Consider n samples $\mathbf{z}_1, \dots, \mathbf{z}_n$, and denote with $\hat{K}_{\mathbf{Z}}$ the corresponding kernel matrix. Suppose that it holds $\|\hat{K}_{\mathbf{Z}}\|_2 \geq n\lambda$ for a constant $\lambda > 0$. Fix $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_l$ samples with*

$$l \geq \frac{8}{3\gamma^2\lambda} \ln \frac{16 \operatorname{tr}_\lambda(\hat{K}_{\mathbf{Z}})}{\gamma}$$

Then, the matrix $ZZ^ - n\lambda I$ is a γ -approximation of $\hat{K}_{\mathbf{Z}} - n\lambda I$ with probability at least $1 - \gamma$, for all $\gamma \in (0, 1)$. Here, $\operatorname{tr}_\lambda(\hat{K}_{\mathbf{Z}})$ is defined as the trace of the matrix $\hat{K}_{\mathbf{Z}}(\hat{K}_{\mathbf{Z}} + n\lambda I)^{-1}$.*

We conclude this section by illustrating the use of random Fourier features to approximate a simple Gaussian kernel. Suppose that we are given a kernel of the form

$$k(\mathbf{z}, \mathbf{z}') := \exp \left\{ -\frac{1}{2} \sigma \|\mathbf{z} - \mathbf{z}'\|_2^2 \right\}.$$

Then, $k(\mathbf{z}, \mathbf{z}')$ can be estimated as in Theorem E.2, with $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_l \sim \mathcal{N}(0, \Sigma)$, with $\Sigma := \sigma^{-1}I$, with I the identity matrix. The functions $\zeta_{k,l}(\mathbf{z})$ can be defined accordingly.

F ADDITIONAL EXPERIMENTS AND SETTINGS

This section contains detailed information on the experiments and additional results.

F.1 DATASET FOR MODEL PERFORMANCE WITH THE USE OF THE HSCIC

The data-generating mechanism corresponding to the results in Fig. 2 is the following:

$$\begin{aligned} \mathbf{Z} &\sim \mathcal{N}(0, 1) & \mathbf{A} &= \mathbf{Z}^2 + \varepsilon_{\mathbf{A}} \\ \mathbf{X} &= \exp \left\{ -\frac{1}{2} \mathbf{A}^2 \right\} \sin(2\mathbf{A}) + 2\mathbf{Z} \frac{1}{5} \varepsilon_{\mathbf{X}} \\ \mathbf{Y} &= \frac{1}{2} \exp \{ -\mathbf{XZ} \} \cdot \sin(2\mathbf{XZ}) + 5\mathbf{A} + \frac{1}{5} \varepsilon_{\mathbf{Y}}, \end{aligned}$$

where $\varepsilon_{\mathbf{A}} \sim \mathcal{N}(0, 1)$ and $\varepsilon_{\mathbf{Y}}, \varepsilon_{\mathbf{X}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 0.1)$.

In the first experiment, Fig. 2 shows the results of feed-forward neural networks consisting of 8 hidden layers with 20 nodes each, connected with a rectified linear activation function (ReLU) and a linear final layer. Mini-batch size of 256 and the Adam optimizer with a learning rate of 10^{-3} for 1000 epochs were used.

F.2 DATASETS AND RESULTS FOR COMPARISON WITH BASELINES

The comparison of our method CIP with the CF1 and CF2 is done on different simulated datasets. These will be referred to as Scenario 1 and Scenario 2. The data generating mechanism corresponding to the results in Fig. 2 (right) is the following:

$$\begin{aligned} \mathbf{Z} &\sim \mathcal{N}(0, 1) & \mathbf{A} &= \exp \left\{ \frac{1}{2} \mathbf{Z}^2 \right\} \cdot \sin(2\mathbf{Z}) + \varepsilon_{\mathbf{A}} \\ \mathbf{X} &= (\mathbf{A} + 0.1\mathbf{Z}) \cdot \varepsilon_{\mathbf{X}} \\ \mathbf{Y} &= \mathbf{A} + \mathbf{X} + 0.1 \cdot \sin(\mathbf{Z}) \end{aligned}$$

Table 1: **Performance of the HSCIC against baselines CF1 and CF2 on two synthetic datasets.** Notably, in both scenarios it is possible to select γ values for which CIP outperforms CF2 in **MSE** and **VCF** simultaneously.

	Scenario 1			Scenario 2		
	MSE $\times 10^6$	HSCIC $\times 10^3$	VCF $\times 10^3$	MSE $\times 10^3$	HSCIC $\times 10^2$	VCF $\times 10^2$
$\gamma = 0.001$	12 \pm 9	45.38 \pm 0.41	54.93 \pm 7.50	0.0006 \pm 0.0002	35.64 \pm 0.32	5.60 \pm 0.03
$\gamma = 0.01$	16 \pm 12	45.35 \pm 0.41	54.57 \pm 7.18	0.0019 \pm 0.0003	35.44 \pm 0.33	5.50 \pm 0.03
$\gamma = 0.1$	32 \pm 20	45.11 \pm 0.43	54.16 \pm 7.58	0.11 \pm 0.006	33.47 \pm 0.36	4.46 \pm 0.04
$\gamma = 0.2$	81 \pm 14	44.78 \pm 0.47	53.59 \pm 7.90	0.42 \pm 0.02	31.38 \pm 0.38	3.52 \pm 0.04
$\gamma = 0.3$	192 \pm 33	43.92 \pm 0.52	52.92 \pm 7.54	0.82 \pm 0.04	29.75 \pm 0.34	2.50 \pm 0.04
$\gamma = 0.4$	384 \pm 58	43.88 \pm 0.57	52.06 \pm 7.25	1.21 \pm 0.05	28.63 \pm 0.33	1.79 \pm 0.03
$\gamma = 0.5$	685 \pm 133	43.26 \pm 0.65	51.64 \pm 7.40	1.56 \pm 0.08	27.81 \pm 0.26	1.1 \pm 0.01
$\gamma = 0.6$	1117 \pm 165	42.47 \pm 0.73	50.96 \pm 7.36	1.84 \pm 0.11	26.87 \pm 0.22	0.79 \pm 0.01
$\gamma = 0.7$	1655 \pm 223	42.11 \pm 0.80	50.31 \pm 7.44	2.11 \pm 0.14	26.08 \pm 0.20	0.49 \pm 0.01
$\gamma = 0.8$	2225 \pm 296	41.87 \pm 0.84	49.76 \pm 7.25	2.37 \pm 0.15	25.27 \pm 0.18	0.31 \pm 0.01
$\gamma = 0.9$	2832 \pm 372	41.52 \pm 0.92	49.17 \pm 7.41	2.58 \pm 0.17	24.64 \pm 0.16	0.21 \pm 0.01
$\gamma = 1.0$	3472 \pm 422	38.37 \pm 0.97	48.71 \pm 7.55	2.77 \pm 0.19	24.21 \pm 0.15	0.14 \pm 0.01
CF1	10321 \pm 72	41.37 \pm 0.58	0 \pm 0.00	4.59 \pm 0.4478	25.01 \pm 0.25	0 \pm 0.00
CF2	2728 \pm 272	41.37 \pm 0.92	59.50 \pm 10.35	3.97 \pm 0.3479	27.03 \pm 0.35	2.62 \pm 0.81

where $\varepsilon_{\mathbf{A}}, \varepsilon_{\mathbf{X}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$ and $\varepsilon_{\mathbf{Y}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 0.1)$. This is referred to as Scenario 1. The data generating mechanism for Scenario 2 is the following:

$$\mathbf{Z} \sim \mathcal{N}(0, 1) \quad \mathbf{A} = \exp\left\{\frac{1}{2}\mathbf{Z}^2\right\} \cdot \sin(2\mathbf{Z}) + \varepsilon_{\mathbf{A}}$$

$$\mathbf{X} = \exp\left\{-\frac{1}{2}\mathbf{A}^2\right\} \cdot \varepsilon_{\mathbf{X}} + 2\mathbf{Z}$$

$$\mathbf{Y} = \frac{1}{2} \sin(\mathbf{Z}\mathbf{X}) \cdot \exp\{-\mathbf{Z}\mathbf{X}\} + \frac{1}{5}\varepsilon_{\mathbf{Y}},$$

where $\varepsilon_{\mathbf{A}}, \varepsilon_{\mathbf{X}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$ and $\varepsilon_{\mathbf{Y}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 0.1)$. Fig. 2 (right) and Table 1 present the average and standard deviation resulting from 9 random seeds runs. For CIP, the same hyperparameters as in the previous setting are used. The MLPs implemented in CF1 and CF2 used for the prediction of $\hat{\mathbf{Y}}$ and the one used for the prediction of the \mathbf{X} residuals in CF2 are all designed with similar architecture and training method. The MLP models consist of 8 hidden layers with 20 nodes each, connected with a rectified linear activation function (ReLU) and a linear final layer. During training, mini-batch size of 64 and the Adam optimizer with a learning rate of 10^{-3} for 200 epochs were used.

F.3 DATASETS AND RESULTS FOR MULTI-DIMENSIONAL VARIABLES EXPERIMENTS

The data-generating mechanisms for the multi-dimensional settings of Fig. 3 are now shown. Given $\dim \mathbf{A} = D_1 \geq 2$, the datasets were generated from:

$$\mathbf{Z} \sim \mathcal{N}(0, 1) \quad \mathbf{A}_i = \mathbf{Z}^2 + \varepsilon_{\mathbf{A}}^i \quad \text{for } i \in \{1, D_1\}$$

$$\mathbf{X} = \exp\left\{-\frac{1}{2}\mathbf{A}_1\right\} + \sum_{i=1}^{D_1} \mathbf{A}_i \cdot \sin(\mathbf{Z}) + 0.1 \cdot \varepsilon_{\mathbf{X}}$$

$$\mathbf{Y} = \exp\left\{-\frac{1}{2}\mathbf{A}_2\right\} \cdot \sum_{i=1}^{D_1} \mathbf{A}_i + \mathbf{X}\mathbf{Z} + 0.1 \cdot \varepsilon_{\mathbf{Y}},$$

where $\varepsilon_{\mathbf{X}}, \varepsilon_{\mathbf{Y}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 0.1)$ and $\varepsilon_{\mathbf{A}}^1, \dots, \varepsilon_{\mathbf{A}}^{D_1} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$. In this experiment, the mini-batch size chosen is 512 and the same hyperparameters are used as in the previous settings. The neural network

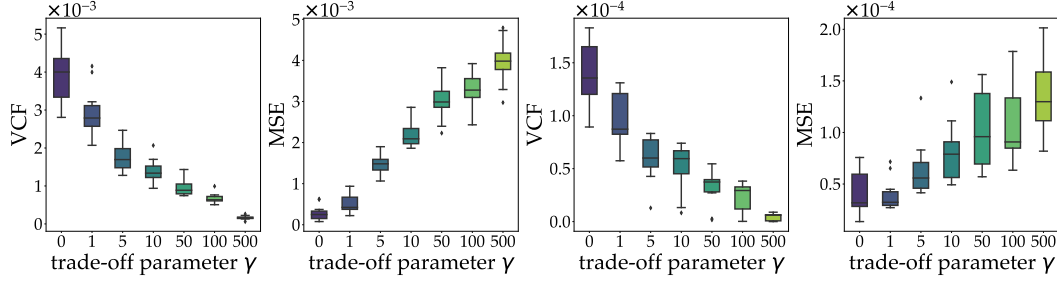


Figure 7: MSE , HSCIC, VCF for increasing dimension of \mathbf{A} on synthetic data from Appendix F.3 with $\text{dimA} = 20$ (left) and $\text{dimA} = 100$ (right). All other variables are one-dimensional.

Table 2: Architecture of the convolutional neural network used for the image dataset, as described in Appendix F.4.

layer	# filters	kernel size	stride size	padding size
convolution	16	5	2	2
max pooling	1	3	2	0
convolution	64	5	1	2
max pooling	1	1	2	0
convolution	64	5	1	2
max pooling	1	2	1	0
convolution	16	5	1	3
max pooling	1	2	2	0

architecture is trained for 800 epochs. Fig. 7 present the results corresponding to 10 random seeds with different values of the trade-off parameter γ corresponding to different values of dimA among $\{15, 100\}$. In all of the box plots, it is evident that there exists a trade-off between the accuracy and counterfactual invariance of the predictor. As the value of γ increases, there is a consistent trend of augmenting counterfactual invariance (as evidenced by the decrease in the VCF metric). Similarly to the previous boxplots visualizations, the boxes represent the interquartile range (IQR), the horizontal line is the median, and whiskers show the minimum and maximum values, excluding the outliers (determined as a function of the inter-quartile range). Outliers are represented in the plot as dots.

F.4 IMAGE DATASET

The simulation procedure for the results shown in Section 4.2 is the following.

$$\begin{aligned}
 \text{shape} &\sim \mathbb{P}(\text{shape}) \\
 \text{y-pos} &\sim \mathbb{P}(\text{y-pos}) \\
 \text{color} &\sim \mathbb{P}(\text{color}) \\
 \text{orientation} &\sim \mathbb{P}(\text{orientation}) \\
 \text{x-pos} &= \text{round}(x), \text{ where } x \sim \mathcal{N}(\text{shape} + \text{y-pos}, 1) \\
 \text{scale} &= \text{round}\left(\left(\frac{\text{x-pos}}{24} + \frac{\text{y-pos}}{24}\right) \cdot \text{shape} + \epsilon_S\right) \\
 \mathbf{Y} &= e^{\text{shape}} \cdot \text{x-pos} + \text{scale}^2 \cdot \sin(\text{y-pos}) + \epsilon_Y,
 \end{aligned}$$

where $\epsilon_S \sim \mathcal{N}(0, 1)$ and $\epsilon_Y \sim \mathcal{N}(0, 0.01)$. The data has been generated via a matching procedure on the original dSprites dataset.

In Table 2, the hyperparameters of the layers of the convolutional neural network are presented. Each of the convolutional groups also has a ReLU activation function and a dropout layer. Two MLP architectures have been used. The former takes as input the observed tabular features. It is composed by two hidden layers of 16 and 8 nodes respectively, connected with ReLU activation functions and dropout layers. The latter takes as input the concatenated outcomes of the CNN and the other MLP. It consists of three hidden layers of 8, 8 and 16 nodes, respectively.

F.5 FAIRNESS WITH CONTINUOUS PROTECTED ATTRIBUTES

The pre-processing of the UCI Adult dataset was based upon the work of [Chiappa & Pacchiano \(2021\)](#). Referring to the causal graph in Fig. 8, a variational autoencoder ([Kingma & Welling, 2014](#)) was trained for each of the unobserved variables \mathbf{H}_m , \mathbf{H}_l and \mathbf{H}_r . The prior distribution of these latent variables is assumed to be standard Gaussian. The posterior distributions $\mathbb{P}(\mathbf{H}_m|V)$, $\mathbb{P}(\mathbf{H}_l|V)$, $\mathbb{P}(\mathbf{H}_r|V)$ are modeled as 10-dimensional Gaussian distributions, whose means and variances are the outputs of the encoder.

The encoder architecture consists of a hidden layer of 20 hidden nodes with hyperbolic tangent activation functions, followed by a linear layer. The decoders have two linear layers with a hyperbolic tangent activation function. The training loss of the variational autoencoder consists of a reconstruction term (Mean-Squared Error for continuous variables and Cross-Entropy Loss for binary ones) and the Kullback–Leibler divergence between the posterior and the prior distribution of the latent variables. For training, we used the Adam optimizer with learning rate of 10^{-2} , 100 epochs, mini-batch size 128.

The predictor $\hat{\mathbf{Y}}$ is the output of a feed-forward neural network consisting of a hidden layer with a hyperbolic tangent activation function and a linear final layer. In the training we used the Adam optimizer with learning rate 10^{-3} , mini-batch size 128, and trained for 100 epochs. The choice of the network architecture is based on the work of [Chiappa & Pacchiano \(2021\)](#).

The estimation of counterfactual outcomes is based on a Monte Carlo approach. Given a data point, 500 values of the unobserved variables are sampled from the estimated posterior distribution. Given an interventional value for A , a counterfactual outcome is estimated for each of the sampled unobserved values. The final counterfactual outcome is estimated as the average of these counterfactual predictions. In this experimental setting, we have $k = 100$ and $d = 1000$.

In the causal graph presented in Fig. 8, \mathbf{A} includes the variables age and gender, \mathbf{C} includes nationality and race, \mathbf{M} marital status, \mathbf{L} level of education, \mathbf{R} the set of the working class, occupation, and hours per week and \mathbf{Y} the income class. Compared to [Chiappa & Pacchiano \(2021\)](#), we include the race variable in the dataset as part of the baseline features \mathbf{C} . The loss function is the same as Eq. (1) but Binary Cross-Entropy loss (\mathcal{L}_{BCE}) is used instead of Mean-Squared Error loss:

$$\mathcal{L}_{\text{CIP}}(\hat{\mathbf{Y}}) = \mathcal{L}_{\text{BCE}}(\hat{\mathbf{Y}}) + \gamma \cdot \text{HSCIC} \left(\hat{\mathbf{Y}}, \{\text{age, gender, marital status, education, work}\} \middle| \mathbf{S} \right), \quad (14)$$

where the set $\mathbf{S} = \{\text{Race, Nationality}\}$ blocks all the non-causal paths from $\mathbf{W} \cup \mathbf{A}$ to \mathbf{Y} . In this example we have $\mathbf{W} = \{\mathbf{C} \cup \mathbf{M} \cup \mathbf{L} \cup \mathbf{R}\}$. The results in Fig. 5 (right) refer to one run with conditioning set $\mathbf{S} = \{\text{Race, Nationality}\}$. The results correspond to 4 random seeds.

F.6 ILLUSTRATING THE CHOICE OF γ

In Section 3.4, we propose to choose γ to obtain a maximal level of CI within a given tolerance on predictive performance. Here, we illustrate results from running the proposed procedure that dynamically selects γ , adjusted to different predefined accuracy thresholds in a classification setting. Specifically, the algorithm chooses the largest γ value that yields an accuracy equal to or better than the threshold. As described the algorithm operates on γ values on a logarithmic scale, thereby

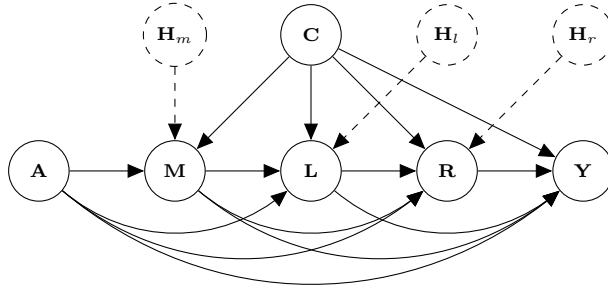


Figure 8: Assumed causal graph for the Adult dataset, as in [Chiappa & Pacchiano \(2021\)](#). The variables \mathbf{H}_m , \mathbf{H}_l , \mathbf{H}_r are unobserved, and jointly trained with the predictor $\hat{\mathbf{Y}}$.

Table 3: Results of MSE and VCF (all times 10^2 for readability) on synthetic data of CIP with trade-off parameters depending on the chosen accuracy threshold.

	VCF $\times 10^2$	HSCIC $\times 10^2$
90% accuracy	3.14 ± 0.92	4.51 ± 0.72
70% accuracy	3.01 ± 0.80	4.44 ± 0.65
1% accuracy	2.91 ± 0.92	4.39 ± 0.42

ensuring a fine-grained search over a wide range of potential trade-off points. Table 3 shows the found trade-offs for tolerated accuracies of 90%, 70%, and 1% in the same setting as Appendix F.1.

G COMPARISON WITH ADDITIONAL BASELINES

In this section, we compare CIP with additional baselines. These include Veitch et al. (2021) and different heuristic methods.

G.1 BASELINE EXPERIMENTS (VEITCH ET AL., 2021)

We provide an experimental comparison against the method by Veitch et al. (2021). To this end, we consider the following data-generating mechanism for the causal structure (see Fig. 1(b)):

$$\begin{aligned}\mathbf{Z} &\sim \mathcal{N}(0, 1) & \mathbf{A} &= \sin(0.1\mathbf{Z}) + \varepsilon_{\mathbf{A}} \\ \mathbf{X} &= \exp\left\{-\frac{1}{2}\mathbf{A}\right\} \sin(\mathbf{A}) + \frac{1}{10}\varepsilon_{\mathbf{X}} \\ \mathbf{Y} &= \frac{1}{10} \exp\{-\mathbf{X}\} \cdot \sin(2\mathbf{XZ}) + \mathbf{A}\mathbf{A} + \frac{1}{10}\varepsilon_{\mathbf{Y}},\end{aligned}$$

where $\varepsilon_{\mathbf{X}}, \varepsilon_{\mathbf{A}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$ and $\varepsilon_{\mathbf{Y}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 0.1)$. The data-generating mechanism of the anti-causal structure is the following (see Fig. 1(c)):

$$\begin{aligned}\mathbf{Z} &\sim \mathcal{N}(0, 1) & \mathbf{A} &= \frac{1}{5} \sin(\mathbf{Z}) + \varepsilon_{\mathbf{A}} \\ \mathbf{Y} &= \frac{1}{10} \sin(\mathbf{Z}) + \varepsilon_{\mathbf{Y}} \\ \mathbf{X} &= \mathbf{A} + \mathbf{Y} + \frac{1}{10}\varepsilon_{\mathbf{X}}\end{aligned}$$

where $\varepsilon_{\mathbf{Y}}, \varepsilon_{\mathbf{A}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 0.1)$ and $\varepsilon_{\mathbf{X}} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$. We compare our method (CIP) against the method by Veitch et al. (2021) using different values for the trade-off parameter γ . In Fig. 1(b-c) the causal and anti-causal graphical settings proposed by Veitch et al. (2021) are presented. In both of these settings there is an unobserved confounder \mathbf{Z} between \mathbf{A} and \mathbf{Y} . The graphical assumptions outlined in Theorem 3.2 of the CIP are not met in the graphical structures under examination, as the confounding path is not effectively blocked by an observed variable (\mathbf{Z} is unobserved). In light of this, it is assumed in our implementation that there is no unobserved confounder. In the graphical structure Fig. 1(b), CIP enforces $\text{HSIC}(\hat{\mathbf{Y}}, \mathbf{A} \cup \mathbf{X})$ to become small, gradually enforcing $\hat{\mathbf{Y}} \perp\!\!\!\perp \mathbf{A} \cup \mathbf{X}$. HSIC is the Hilbert-Schmidt Independence Criterion, which is commonly used to promote independence (see, i.e., Gretton et al. (2005); Fukumizu et al. (2007)). Veitch et al. (2021) enforces as independence criterion $\text{HSIC}(\hat{\mathbf{Y}}, \mathbf{A})$, which is implied by the independence enforced in CIP. In the anti-causal graphical setting presented in Fig. 1(c), the objective term used in CIP is $\text{HSCIC}(\hat{\mathbf{Y}}, \mathbf{A} \mid \mathbf{X})$, while in the method of Veitch et al. (2021) is $\text{HSCIC}(\hat{\mathbf{Y}}, \mathbf{A} \mid \mathbf{Y})$. In Table 4, the results of accuracy and VCF are presented.

In the experiments, the predictor $\hat{\mathbf{Y}}$ is a feed-forward neural network consisting of 8 hidden layers with 20 nodes each, connected with a rectified linear activation function (ReLU) and a linear final layer. Mini-batch size of 256 and the Adam optimizer with a learning rate of 10^{-4} for 500 epochs were used.

Table 4: **Results of the MSE, VCF of CIP and the baseline (Veitch et al., 2021)** applied to the causal and anti-causal structure in Fig. 1(b-c). Although the graphical assumptions are not satisfied, CIP shows an overall decrease of VCF in both of the graphical structures, performing on par with the baseline Veitch et al. (2021) in terms of accuracy and counterfactual invariance.

	CIP		Veitch et al. (2021)	
	MSE $\times 10^2$	VCF	MSE $\times 10^2$	VCF
$\gamma = 0.5$	4.58 ± 0.31	0.19 ± 0.02	4.50 ± 0.40	0.19 ± 0.02
$\gamma = 1.0$	5.60 ± 0.36	0.18 ± 0.01	5.45 ± 0.41	0.18 ± 0.02

	CIP		Veitch et al. (2021)	
	MSE $\times 10^2$	VCF	MSE $\times 10^2$	VCF
$\gamma = 0.5$	1.16 ± 0.01	1.69 ± 0.16	1.01 ± 0.01	1.71 ± 0.26
$\gamma = 1.0$	1.37 ± 0.02	1.48 ± 0.19	0.99 ± 0.01	1.88 ± 0.28

Table 5: Results of MSE and VCF (all times 10^2 for readability) on synthetic data of CIP with trade-off parameters $\gamma = 0.5$ and $\gamma = 1$ with the heuristic methods *data augmentation* and *causal-based data augmentation* and *naive prediction*.

	VCF $\times 10^3$	MSE $\times 10^3$
data augmentation	3.12 ± 0.16	0.03 ± 0.01
causal-based data augmentation	3.04 ± 0.16	0.13 ± 0.12
CIP ($\gamma = 0.5$)	1.05 ± 0.13	1.64 ± 0.22
CIP ($\gamma = 1.0$)	0.35 ± 0.19	2.50 ± 0.72
naive prediction (ignore A)	9.01 ± 0.02	3.01 ± 0.91

G.2 COMPARISON BASELINES HEURISTIC METHODS

We provide an experimental comparison of the proposed method (CIP) with some heuristic methods, specifically data-augmentation-based methods. We consider the same data-generating procedure and causal structure as presented in Appendix F.1. The heuristic methods considered are *data augmentation* and *causal-based data augmentation*. In the former, data augmentation is performed by generating $N = 50$ samples for every data-point by sampling new values of \mathbf{A} as $a_1, \dots, a_N \stackrel{i.i.d}{\sim} \mathbb{P}_{\mathbf{A}}$ and leaving $\mathbf{Z}, \mathbf{X}, \mathbf{Y}$ **unchanged**. Differently, in the latter *causal-based data augmentation* method, we also take into account the causal structure given by the known DAG. Indeed, when manipulating the variable \mathbf{A} , its descendants (in this example \mathbf{X}) will also change. In this experiment, a predictor for \mathbf{X} as $\hat{\mathbf{X}} = f_{\theta}(\mathbf{A}, \mathbf{Z})$ is trained on 80% of the original dataset. In the data augmentation mechanism, for every data-point $\{a, x, z, y\}$, $N = 50$ samples are generated by sampling new values of \mathbf{A} as $a_1, \dots, a_N \stackrel{i.i.d}{\sim} \mathbb{P}_{\mathbf{A}}$, estimating the values of \mathbf{X} as $x_1 = f_{\theta}(a_1, z), \dots, x_N = f_{\theta}(a_N, z)$, while leaving the values of \mathbf{Z} and \mathbf{Y} unchanged. Heuristic methods such as data-augmentation methods do not theoretically guarantee to provide counterfactually invariant predictors. The results of an empirical comparison are shown in Table 5 with the average and standard deviations after 5 random seeds. It can be shown that these theoretical insights are supported by experimental results, as the VCF metric measure counterfactual invariance is lower in both of the two settings of the CIP ($\gamma = \frac{1}{2}$ and $\gamma = 1$).

A dataset of $n = 3000$ is used, along with $k = 500$ and $d = 500$. The architecture for predicting \mathbf{X} and \mathbf{Y} are feed-forward neural networks consisting of 8 hidden layers with 20 nodes each, connected with a rectified linear activation function (ReLU) and linear final layer. Mini-batch size of 256 and the Adam optimizer with a learning rate of 10^{-3} for 100 epochs were used.