

# SYMMETRY-CONSTRAINED CAUSAL PARTIAL IDENTIFICATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We present a novel framework for using knowledge of data symmetries to sharpen bounds in causal *partial identification* (PI). The causal effect of the treatment  $X$  on outcome  $Y$  is generally not identifiable from observational data alone if their common causes, also known as confounders, are unobserved. PI entails estimating bounds on such treatment effects by solving a constrained optimization problem that encodes different assumptions imposed on data generation. PI has use in many application domains where such bounds are sufficient to inform policy decisions, even if the treatment effect itself is not identifiable. We show that knowledge of symmetries in data generation—formalized as invariance under transformation groups—provides additional constraints that tighten these bounds. We operationalize this insight through two approaches: (i) adding explicit invariance error constraints to existing PI methods, and (ii) applying symmetry-preserving *data augmentation* (DA) as a pre-processing step. Under a linear Gaussian model, we show that the later yields bounds that provably valid (containing the true causal effect), sharper (smaller identified sets), and more robust (lower worst-case error). The key mechanism being that randomized symmetry transformations introduce exogenous variation in  $X$  that cannot be attributed to confounding, thereby reducing ambiguity in the identified set. Experiments on synthetic and real data validate our approach. More broadly, our findings establish known data symmetries—ubiquitously employed in DA for variance reduction—can be repurposed as a principled tool for causal inference when point-identification is impossible.

## 1 INTRODUCTION

The problem of regression in machine learning aims to fit a model to observational  $(X, Y)$  data that predicts outcome  $Y$  from treatment  $X$ . Improving the generalization of such predictors to unlabeled samples of  $X$  often requires regularization techniques like *data augmentation* (DA) (Vapnik, 1998; Shorten & Khoshgoftaar, 2019; Lyle et al., 2020). However, such predictive models are generally not causal: the statistical relationship between  $X$  and  $Y$  may be driven by unobserved common causes, i.e. confounders, rather than the true influence of  $X$  on  $Y$ . The gold standard for eliminating confounding is direct intervention, i.e. explicit randomization of  $X$  during data generation (Peters et al., 2017; Pearl, 2009). Since these are often inaccessible, a common workaround is to correct for confounding via auxiliary variable (Zhang et al., 2023). However, these too may be insufficient to recover the causal effect (Kilbertus et al., 2020), or scarce in many applications (Akbar et al., 2025).

In such cases, identifying the true causal effect is not possible from observational data alone. *Partial identification* (PI) offers a principled alternative by computing bounds guaranteed to contain the true effect (Padh et al., 2023)—often sufficient for decision-making even without point-identification.. These bounds are obtained by solving an optimization problem whose constraints encode assumptions about how the data were generated. The informativeness of the bounds depends entirely on the strength and structure of these constraints.

This paper introduces known data symmetries—formalized as invariance under transformations of  $X$ —as a new source of constraints for PI. Such symmetries are ubiquitous in scientific and causal modeling, ranging from geometric stabilities in physical systems (Bronstein et al., 2021; Satorras et al., 2022) to semantic invariance in natural language (Veitch et al., 2021a) and permutation in-

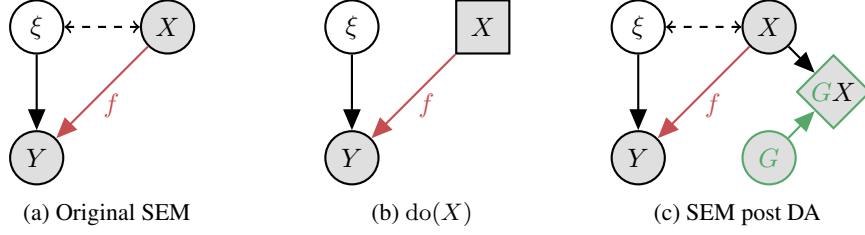


Figure 1: Graphs of respective SEMs. (a) The original SEM from Eq. (1) with confounded  $(X, Y)$ . (b) The original SEM post intervention on  $X$ . (c) The original SEM post DA application.

variance in exchangeable data (Veitch et al., 2021b). We demonstrate that enforcing this structural knowledge restricts the set of plausible causal hypotheses, effectively pruning the identified set. We operationalize this insight in two complementary ways:

1. **Explicit constraints via invariance error (Section 4.1):** we propose bounding the invariance error under specified transformations as an explicit constraint for improving PI.
2. **Implicit constraints via DA (Section 4.2):** we show the effectiveness of a simple symmetry-preserving DA based pre-processing step before running off-the-shelf PI solvers.

The first approach gives strict improvement in PI by construction, is amenable to modern Monte Carlo, gradient-based solvers, is not restricted to any specific hypothesis class and allows tunability to easily handle symmetry miss-specification. The later provides a cheap tool that can be composed with black-box PI solvers, offering guarantees under the linear Gaussian regime with well-specified symmetries. Our methods are simple to apply, compatible with existing solvers, and re-purpose the rather pervasive ML tool of data symmetries as a practical approach for strengthening causal conclusions when point-identification not possible.

## 2 PRELIMINARIES

### 2.1 STATISTICAL VS. CAUSAL INFERENCE

Consider random treatment  $X$ , outcome  $Y$  taking values in  $\mathcal{X} \subseteq \mathbb{R}^m$ ,  $\mathcal{Y} \subseteq \mathbb{R}$  respectively. The function  $f \in \mathcal{H} := \{h : \mathcal{X} \rightarrow \mathcal{Y}\}$  defines their causal relationship via a *structural equation model (SEM)*

$$Y = f(X) + \xi, \quad \mathbb{E}[\xi] = 0. \quad (1)$$

We want to estimate  $f$  given a dataset  $\mathcal{D} := \{(x_i, y_i)\}_{i=0}^n$  of  $n$  samples from the distribution  $\mathbb{P}_{X,Y}$ .

With the assumption  $X \perp\!\!\!\perp \xi$ , we have  $\mathbb{E}[Y|X = \mathbf{x}] = f(\mathbf{x})$  in Eq. (1). *Statistical inference* entails identifying precisely the Bayes optimal predictor  $\mathbb{E}[Y|X = \mathbf{x}]$  from  $\mathcal{D}$  by minimizing an empirical version of the *statistical risk* over hypotheses  $h \in \mathcal{H}$  for some proper, convex loss  $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+$ ,

$$R_{\text{erm}}(h) := \mathbb{E}[\ell(Y, h(X))]. \quad (2)$$

Then, for a sufficiently rich hypothesis class, the minimizer  $h_{\text{erm}}$  gives an unbiased estimation of  $f$ .

However, the residual  $\xi$  in Eq. (1) may generally be correlated with  $X$ , i.e.,  $\mathbb{E}[\xi|X] \neq 0$ , so that the conditional  $\mathbb{E}[Y|X = \mathbf{x}]$  now gives a biased estimate of  $f(\mathbf{x})$  (Pearl, 2009; Peters et al., 2017). This correlation arises due to unobserved common causes of  $X$  and  $Y$ , known as *confounders*. We say that  $X$  and  $Y$  are confounded and refer to the resulting bias as the *confounding bias* (Pearl, 2009). *Causal inference* entails adjusting for this bias to identify  $f$ , or at the very least account for it by finding bounds on  $f$  should identification not be possible. Both approaches are outlined below.

### 2.2 INTERVENTION FOR CAUSAL EFFECT IDENTIFICATION

We can make  $X$  and the residual  $\xi$  uncorrelated via an *intervention*  $\text{do}(X := X')$  that explicitly sets  $X$  to some independently sampled  $X'$  in Eq. (1) during data generation. The induced distribution, referred to as the *interventional distribution*, is represented by  $\mathbb{P}_{X,Y}^{\text{do}(X:=X')}$ . We use the shorthand

notation  $\text{do}(X)$  for an intervention where  $X' \sim \mathbb{P}_X$ , under which the objective from Eq. (2) now defines the *causal risk* (Kania & Wit, 2023; Vankadara et al., 2022; Janzing & Schölkopf, 2018b) as

$$R_{\text{erm}}^{\text{do}(X)}(h) := \mathbb{E}^{\text{do}(X)}[\ell(Y, h(X))]. \quad (3)$$

The target estimand of Eq. (3) is the *average treatment effect (ATE)*  $\mathbb{E}^{\text{do}(X:=\mathbf{x})}[Y | X = \mathbf{x}]$  which equals  $f(\mathbf{x})$  for the SEM under consideration in Eq. (1). Minimizers of Eq. (3) therefore give an unbiased estimation of  $f$ . To better capture the *estimation error* for a candidate hypothesis  $h \in \mathcal{H}$ , we use the *causal excess risk* (Vankadara et al., 2022) by removing irreducible noise from Eq. (3) as

$$E^{\text{do}(X)}(h) := R_{\text{erm}}^{\text{do}(X)}(h) - R_{\text{erm}}^{\text{do}(X)}(f).$$

Since interventions are often inaccessible for computing the risk Eq. (3), estimating  $f$  usually relies on access to the full joint distribution  $\mathbb{P}$  of  $(X, Y, \xi)$  via *back-door adjustment* (Xu & Gretton, 2022)

$$h_{\text{adj}}^{\mathbb{P}}(\mathbf{x}) := \mathbb{E}_{\xi}[\mathbb{E}[Y | X = \mathbf{x}, \xi]], \quad (X, Y, \xi) \sim \mathbb{P}.$$

### 2.3 PARTIAL IDENTIFICATION AND SENSITIVITY ANALYSIS

For unobserved noise  $\xi$ , identification of  $f$  is generally not possible from  $\mathbb{P}_{X,Y}$  alone. Nevertheless, given assumptions on the data generating process in Eq. (1), we can do *partial identification (PI)* (Padh et al., 2023) of  $f$  by considering all the joint distributions  $\mathbb{Q}$  consistent with said assumptions,

$$\mathcal{Q}_{\text{pi}}(\mathbb{P}_{X,Y}) := \left\{ \mathbb{Q} \in \mathcal{C}_{\text{pi}} \mid \mathbb{Q}_{X,Y} = \mathbb{P}_{X,Y} \right\},$$

where the constraint set  $\mathcal{C}_{\text{pi}}$  encodes our assumptions. If correctly specified,  $\mathbb{P} \in \mathcal{C}_{\text{pi}}$  and the following set  $\mathcal{H}_{\text{pi}}$  of candidate hypotheses contains the true solution  $f$ , or the interval  $\mathcal{H}_{\text{pi}}(\mathbf{x})$  holds  $f(\mathbf{x})$ ,

$$\mathcal{H}_{\text{pi}} := \left\{ h_{\text{adj}}^{\mathbb{Q}} \mid \mathbb{Q} \in \mathcal{Q}_{\text{pi}} \right\}, \quad \mathcal{H}_{\text{pi}}(\mathbf{x}) := \left\{ h_{\text{adj}}^{\mathbb{Q}}(\mathbf{x}) \mid \mathbb{Q} \in \mathcal{Q}_{\text{pi}} \right\},$$

where  $\mathcal{Q}_{\text{pi}}$  is shorthand for  $\mathcal{Q}_{\text{pi}}(\mathbb{P}_{X,Y})$ . Computing the interval  $\mathcal{H}_{\text{pi}}(\mathbf{x})$  at  $\mathbf{x}$  is often more practical than characterizing the set  $\mathcal{H}_{\text{pi}}$ , since it amounts to solving two constrained optimization problems as

$$\mathcal{H}_{\text{pi}}(\mathbf{x}) = \left[ \min_{\mathbb{Q} \in \mathcal{Q}_{\text{pi}}} h_{\text{adj}}^{\mathbb{Q}}(\mathbf{x}), \max_{\mathbb{Q} \in \mathcal{Q}_{\text{pi}}} h_{\text{adj}}^{\mathbb{Q}}(\mathbf{x}) \right]. \quad (4)$$

In either case, we want the identified sets to (i) contain the true solution, (ii) be as small as possible.

The constraint set may also be parameterized as  $\mathcal{C}_{\text{pi}}(\Gamma)$  to conduct *sensitivity analyses* (Frauen et al., 2024) by varying parameters  $\Gamma$  to see how  $\mathcal{H}_{\text{pi}}$ ,  $\mathcal{H}_{\text{pi}}(\mathbf{x})$  evolve as assumptions are relaxed/tightened.

Lastly, since we are now discussing hypothesis sets, we define the two appropriate evaluation metrics

$$E_{\text{approx}}^{\text{do}(X)}(\mathcal{Q}_{\text{pi}}) := \min_{\mathbb{Q} \in \mathcal{Q}_{\text{pi}}} E^{\text{do}(X)}(h_{\text{adj}}^{\mathbb{Q}}), \quad E_{\text{worst}}^{\text{do}(X)}(\mathcal{Q}_{\text{pi}}) := \max_{\mathbb{Q} \in \mathcal{Q}_{\text{pi}}} E^{\text{do}(X)}(h_{\text{adj}}^{\mathbb{Q}}).$$

The *approximation error*  $E_{\text{approx}}^{\text{do}(X)}(\mathcal{Q}_{\text{pi}})$  measures how far the target  $f$  is from  $\mathcal{H}_{\text{pi}}$  (Brown & Ali, 2024), and the *worst-case excess risk*  $E_{\text{worst}}^{\text{do}(X)}(\mathcal{Q}_{\text{pi}})$  upper bounds the performance of the identified set  $\mathcal{H}_{\text{pi}}$  relative to the target  $f$ . Similarly,  $\mathcal{H}_{\text{pi}}(\mathbf{x})$  is evaluated using  $E_{\text{approx}}^{\text{do}(\mathbf{x})}(\mathcal{Q}_{\text{pi}})$  and  $E_{\text{worst}}^{\text{do}(\mathbf{x})}(\mathcal{Q}_{\text{pi}})$ .

**Choice of constraints.** The nature and construction of  $\mathcal{C}_{\text{pi}}$  often depends on domain knowledge. Popular approaches involve bounding the spurious correlation between  $X, Y$ , including the sensitivity model by Rosenbaum (2002) which parameterizes the strength of unmeasured confounding through odds ratios, its generalization of the Marginal Sensitivity Model (MSM) by Tan (2006) that does the same using propensity scores and the partial R-squared approach by Cinelli & Hazlett (2020) bounds the proportion of variance explained by unobserved confounders. More recently, Fan et al. (2024); Guo et al. (2022) formulated the PI problem as *robust optimization (RO)* over  $\mathcal{Q}_{\text{pi}}$  constructed as a total variation ball around the observational distribution  $\mathbb{P}_{X,Y}$ , and Meresht et al. (2022) similarly uses Wasserstein constraints. An equivalent approach to modeling the confounding is to instead model the random function  $f_{\xi}(\cdot) := f(\cdot) + \xi$  itself, also known as the *response function* (Padh et al., 2023). Hu et al. (2021) modeled these using generative adversarial networks (GANs) to then match  $\mathbb{P}_{X,Y}$  in distribution. Most of these methods can also leverage auxiliary variables in addition to  $X, Y$  for imposing constraints in the form of conditional independences to sharpen bounds. Of note is the instrumental variable (IV) based PI by Balke & Pearl (1997) for when  $\xi$  arbitrarily influences  $Y$  instead of the additive model in Eq. (1). Modern neural-network based variants for continuous, high-dimensional treatments and/or IVs are explored by Schweisthal et al. (2025); Kilbertus et al. (2020); Hu et al. (2021); Padh et al. (2023); Meresht et al. (2022); Gunsilius (2020).

## 2.4 DATA SYMMETRIES AND INVARIANCE

For finite samples, the technique of *data augmentation (DA)* is used to reduce estimation variance (Lyle et al., 2020; Chen et al., 2020) in statistical inference. This is achieved by applying random transformations  $G \sim \mathbb{P}_G$  to the data, generating multiple transformed samples  $(G\mathbf{x}_i, y_i)$  from each original sample  $(\mathbf{x}_i, y_i) \in \mathcal{D}$ , thereby increasing variability in the data for statistical risk evaluation,

$$R_{\text{da+erm}}(h) := \mathbb{E}[\ell(Y, h(GX))]. \quad (5)$$

In this work we restrict ourselves to DA with respect to which  $f$  is invariant (Lyle et al., 2020; Chen et al., 2020). The action of a group  $\mathcal{G}$  is a mapping  $\alpha : \mathcal{X} \times \mathcal{G} \rightarrow \mathcal{X}$  compatible with the group operation. Writing  $g\mathbf{x} := \alpha(\mathbf{x}, g)$  as shorthand, we say that  $f$  is *invariant* under  $\mathcal{G}$  (or  $\mathcal{G}$ -invariant) if

$$f(g\mathbf{x}) = f(\mathbf{x}), \quad \forall (g, \mathbf{x}) \in \mathcal{G} \times \mathcal{X}.$$

Group  $\mathcal{G}$  has a (unique) normalized Haar measure,  $\mathbb{P}_G$  the corresponding distribution defined over it.

Of course one needs to have prior knowledge about the symmetries of  $f$  to construct such a DA. We argue that the popularity of this modeling assumption in the DA and invariance literature (Lyle et al., 2020; Chen et al., 2020) is precisely because such symmetries are already established in many application domains. For example, when classifying images of cats and dogs we already know that whatever the true labeling function may be, it would certainly be invariant to rotations on the images.  $G$  would then represent the random rotation angle, whereas  $G\mathbf{x}$  would be the rotated image  $\mathbf{x}$ .

While DA is canonically used to mitigate finite-sample estimation variance, our focus is primarily on the infinite-sample setting, and we present Eq. (5) and subsequent theoretical results in that context. Nonetheless, increasing sample size via DA also bears on our work, a point we shall briefly discuss. Section 4.2 also makes the following assumption which characterizes many standard DA operations.

**Assumption 1** (unbiased group action). The group action  $G \sim \mathbb{P}_G$  is identity-centered, meaning

$$\mathbb{E}[G\mathbf{x}] = \mathbf{x}, \quad \forall \mathbf{x} \in \mathcal{X}.$$

**Lemma 1** (added exogenous variation with DA). *Under Assumption 1,  $G$  inflates the data variance,*

$$\Delta := \Sigma_{GX} - \Sigma_X \succcurlyeq \mathbf{0}, \quad \text{equality iff } GX = X \text{ a.s.}$$

*Proof.* See Appendix A.5 for the proof.  $\square$

## 3 CAUSAL IMPLICATIONS OF DATA SYMMETRIES

Crucially, the random group action  $G$  from Lemma 1 introduces additional *exogenous* variation in  $X$  that is independent of other system variables. Consequently, Akbar et al. (2025) showed that for a  $\mathcal{G}$ -invariant target function  $f$ , the transformation  $GX$  simulates a *soft* intervention on  $X$ —perturbing  $X$  to weaken the confounding association  $X \leftrightarrow \xi$  while preserving the causal mechanism  $X \rightarrow Y$ . To formalize how this improves point estimation of  $f$ , Akbar et al. (2025) employs the following linear version of Eq. (1), which also serves as the basis for our subsequent PI analysis in Section 4.2.

**Assumption 2** (a linear, Gaussian SEM). The SEM Eq. (1) is centered, joint Gaussian with  $\mathbf{f} \in \mathbb{R}^m$ ,

$$Y = \mathbf{f}^\top X + \xi.$$

In this setting, the causal estimation error (excess risk) under a squared loss takes the following form:

$$E^{\text{do}(X)}(\mathbf{h}) = \|\mathbf{h} - \mathbf{f}\|_{\Sigma_X}^2, \quad E^{\text{do}(\mathbf{x})}(\mathbf{h}) = (\mathbf{h}^\top \mathbf{x} - \mathbf{f}^\top \mathbf{x})^2, \quad (6)$$

where we overload the notation  $\text{do}(\mathbf{x})$  (as opposed to  $\text{do}(X)$ ) to mean the *hard* intervention  $\text{do}(X := \mathbf{x})$ , i.e. fixing  $X$  to a constant  $\mathbf{x}$  during data generation. Similar formulations have been used to measure causal error (Vankadara et al., 2022; Kania & Wit, 2023; Akbar et al., 2025) or quantify confounding strength (Janzing, 2019; Janzing & Schölkopf, 2018a;b). (Akbar et al., 2025) established the following result, which directly bears on our work:

**Proposition 1** (estimation with DA (Akbar et al. (2025) lifted)). *For  $\mathcal{G}$ -inv.  $f$ , Assumptions 1 and 2,*

$$0 \leq \frac{\kappa}{1+\kappa} \cdot \underbrace{\|\Pi_{\Delta}(\mathbf{h}_{\text{erm}} - \mathbf{f})\|_{\Sigma_X}^2}_{\text{estimation error within range}(\Delta)} \leq E^{\text{do}(X)}(\mathbf{h}_{\text{erm}}) - E^{\text{do}(X)}(\mathbf{h}_{\text{da+erm}}),$$

$$\leq \|\Pi_{\Delta}(\mathbf{h}_{\text{erm}} - \mathbf{f})\|_{\Sigma_X}^2, \quad \text{eq. iff} \quad \underbrace{\Delta \perp \Sigma_{X,\xi}}_{\text{DA orthogonal to confounding}},$$

where  $\kappa := \lambda_{\min}^+(\Sigma_X^{-1}\Delta) < \infty$  represents the lowest positive eigenvalue of the product  $\Sigma_X^{-1}\Delta$ .

*Proof.* See Appendix A.3, cf. (Akbar et al., 2025, Thm. 3) for the proof.  $\square$

Essentially, for  $\mathcal{G}$ -invariant  $f$ , DA *dominates* ERM on causal estimation—performing strictly better iff it aligns with the confounding within  $X$ , but never worse. Note that in Proposition 1, for the “large DA” regime, which we define as  $\kappa \rightarrow \infty$ , the lower-bound approaches the upper-bound, which is simply the sq.-norm of the *projection*  $\Pi_{\Delta}(\cdot)$  of estimation bias  $(\mathbf{h}_{\text{erm}} - \mathbf{f})$  onto  $\text{range}(\Delta)$ . This confirms that identification is generally not possible in this setting. Therefore the principled approach is to undertake PI of  $f$  instead of the point-estimation approach by Akbar et al. (2025).

This motivates our current work, where we leverage knowledge of symmetries in data generation to improve partial identification and/or sensitivity analysis of  $f$ , as discussed in the following section.

## 4 SYMMETRY-CONSTRAINED PARTIAL IDENTIFICATION

Our objective is to leverage symmetry knowledge to restrict the identified sets  $\mathcal{H}_{\text{pi}}$  and  $\mathcal{H}_{\text{pi}}(\mathbf{x})$ . We give two strategies to operationalize this: (i) integrating an explicit invariance error constraint into the optimization, and (ii) inducing implicit regularization through data augmentation pre-processing.

### 4.1 EXPLICIT CONSTRAINT WITH INVARIANCE ERROR

We start off by considering the most obvious approach to incorporate symmetry knowledge into PI—add an explicit invariance error constraint to any baseline PI method defined by a constraint set  $\mathcal{C}_{\text{pi}}$ ,

$$E_{\text{inv}}(h) := \mathbb{E} \left[ \|h(X) - h(GX)\|^2 \right],$$

$$\mathcal{Q}_{\text{inv+pi}}(\mathbb{P}_{X,Y}) := \left\{ \mathbb{Q} \in \mathcal{C}_{\text{pi}} \mid \mathbb{Q}_{X,Y} = \mathbb{P}_{X,Y}, \quad E_{\text{inv}}\left(h_{\text{adj}}^{\mathbb{Q}}\right) \leq \epsilon \right\}.$$

**Remark 1** (sharper, robust bounds with invariance error). By construction, subset inclusion holds:

$$\mathcal{H}_{\text{inv+pi}} \subseteq \mathcal{H}_{\text{pi}}, \quad \mathcal{H}_{\text{inv+pi}}(\mathbf{x}) \subseteq \mathcal{H}_{\text{pi}}(\mathbf{x}).$$

Consequently, this guarantees that the volume of the identified set and the corresponding worst-case excess risk does not increase. Note that due to this same set inclusion, the approximation error generally cannot decrease, and may even potentially increase if  $\mathcal{C}_{\text{pi}}$  does not contain the true distribution  $\mathbb{P}$ . For  $\epsilon = 0$ , these metrics are equal to “large DA” regime results in Sections 3 and 4.2.

Nevertheless, whenever the baseline PI constraints  $\mathcal{C}_{\text{pi}}$  are valid and  $E_{\text{inv}}(f) \leq \epsilon$  holds,  $\mathcal{Q}_{\text{inv+pi}}$  guarantees validity. Furthermore, the parameter  $\epsilon$  enables sensitivity analysis, allowing us to inspect how the bounds evolve as we vary the assumed invariance error in our choice of transformations  $G$ . Of course we can similarly use other formulations for  $E_{\text{inv}}$ , such as ones stated in Yang et al. (2019), or restrict ourselves to a hypothesis class that follows our symmetry by design (Cohen & Welling, 2016). However, the later may be restrictive in terms of compatibility with standard PI methods.

While our experiments discuss settings where Eq. (4) for  $\mathcal{H}_{\text{inv+pi}}(\mathbf{x})$  can be solved via convex programming, we emphasize that this explicit constraint formulation is fully compatible with modern deep learning-based PI. Since the functional  $E_{\text{inv}}$  is differentiable and amenable to Monte Carlo evaluation, it can be readily incorporated as a regularizer in augmented Lagrangian and/or gradient-based solvers (Padh et al., 2023; Kilbertus et al., 2020; Meresht et al., 2022; Hu et al., 2021).

Despite this compatibility, incorporating the invariance error constraint still requires modifying the solver logic or objective function. This imposes an implementation burden and precludes the use of specialized or “black-box” PI software where the internal optimization is fixed. This limitation motivates our second approach—a simple data pre-processing strategy that implicitly impose symmetry

constraints by simply feeding augmented data into standard off-the-shelf PI methods. Furthermore, when modeling complex, high-dimensional data, enforcing non-convex invariance constraints during optimization is often more expensive and notoriously unstable (Schweisthal et al., 2025; Padh et al., 2023) as opposed to a simple data pre-processing step.

#### 4.2 IMPLICIT CONSTRAINT WITH DA PRE-PROCESSING

We draw inspiration from IV methods, where “strong” instruments—those inducing significant exogenous variation in  $X$ —are known to yield sharper identification bounds compared to weak instruments (Kilbertus et al., 2020; Padh et al., 2023). This motivates our central inquiry in this section:

*Can the synthetic exogenous variation introduced by DA similarly sharpen PI?*

As with Akbar et al. (2025), which we extend now to the PI setting, the fundamental mechanism for PI sharpening is Lemma 1. Our main insight into why DA aids PI is summarized as follows:

- (i) **Statistical Efficiency:** Most straightforwardly, DA grows effective data size, quelling sampling variation and finite-sample errors; key sources of uncertainty in PI (Imbens & Manski, 2004).
- (ii) **Sharper Bounds:** DA adds variation in  $X$  that is explicitly exogenous, and therefore cannot be attributed to confounding. This reduces ambiguity in PI, which leads to sharper bounds.
- (iii) **Robust Bounds:** By perturbing spurious features, DA reduces confounding bias, centering and contracting the PI bounds around the true solution. This directly minimizes the worst-case error.
- (iv) **Valid Bounds:** Crucially,  $\mathcal{G}$ -invariance of  $f$  guarantees valid bounds with DA if  $\mathcal{C}_{\text{pi}}$  is valid.

We elaborate these via analysis of the linear model from Assumption 2. But first we explicitly define the composition  $\mathcal{Q}_{\text{da+pi}}$  of DA and PI, as well as the specific PI model that we use for our analysis,

$$\mathcal{Q}_{\text{da+pi}}(\mathbb{P}_{X,Y}) := \mathcal{Q}_{\text{pi}}(\mathbb{P}_{GX,Y}).$$

**Assumption 3** (a bounded confounding sensitivity model). Consider the following constraint set.

$$\mathcal{C}_{\text{pi}}(\Gamma) := \left\{ \mathbb{Q} = \mathcal{N}(\mathbf{0}, \cdot) \mid \frac{\text{Var}(\mathbb{E}[\xi | X])}{\text{Var}(\xi)} \leq \Gamma, \quad \text{Var}(\xi) \leq \Gamma_0 \right\}, \quad \Gamma := [\Gamma_0, \Gamma]^\top,$$

where *confounding strength*  $\Gamma \geq 0$  determines our assumption on the variation in  $\xi$  explained by  $X$ .

Assumption 3 adopts the widely used partial R-squared sensitivity model Cinelli & Hazlett (2019), itself a generalization of the classic Rosenbaum (2002). While we employ this model in our analyses, we do not necessarily restrict ourselves to it—under the linear Gaussian setting of Assumption 2, several families of PI and sensitivity models reduce to ellipsoidal constraints equivalent to the form:

**Lemma 2** (characterizing the identified set in a linear, Gaussian case). *Under Assumptions 2 and 3,*

$$\mathcal{H}_{\text{pi}} = \left\{ \mathbf{h} \mid \|\mathbf{h} - \mathbf{h}_{\text{erm}}\|_{\Sigma_X}^2 \leq r(\Gamma)^2 \right\},$$

where the ellipsoid radius  $r(\Gamma) \geq 0$  depends on the choice of constraint parameters. Furthermore,

$$\mathcal{H}_{\text{pi}}(\mathbf{x}) = \left[ \mathbf{h}_{\text{erm}}^\top \mathbf{x} - r(\Gamma) \cdot \|\mathbf{x}\|_{\Sigma_X^{-1}}, \quad \mathbf{h}_{\text{erm}}^\top \mathbf{x} + r(\Gamma) \cdot \|\mathbf{x}\|_{\Sigma_X^{-1}} \right].$$

*Proof.* See Appendix A.5 for the proof.  $\square$

Our results thus carry broader implications for PI and sensitivity analysis, as we discuss in Section 7.

##### 4.2.1 BETTER BOUNDS WITH DATA AUGMENTATION

First and foremost, we investigate if the post-DA bounds are, in some way, better than the baseline PI bounds. That is, if this exercise is useful at all. We present two results to support this claim.

**Proposition 2** (sharper bounds with DA). *For Assumptions 1 to 3, Lebesgue measure (volume)  $|\cdot|$ ,*

$$\frac{|\mathcal{H}_{\text{da+pi}}|}{|\mathcal{H}_{\text{pi}}|} = \sqrt{\frac{\det \Sigma_X}{\det \Sigma_{GX}}} < 1, \quad \frac{|\mathcal{H}_{\text{da+pi}}(\mathbf{x})|}{|\mathcal{H}_{\text{pi}}(\mathbf{x})|} = \frac{\|\mathbf{x}\|_{\Sigma_{GX}^{-1}}}{\|\mathbf{x}\|_{\Sigma_X^{-1}}} \leq 1, \quad \text{equality iff} \quad \mathbf{x} \perp \Delta.$$

*Proof.* See Appendix A.4 for the proof.  $\square$



Proposition 2 states that the hypothesis set  $\mathcal{H}_{\text{da+pi}}$  is strictly smaller than the baseline  $\mathcal{H}_{\text{pi}}$ . The same holds true for the intervals  $\mathcal{H}_{\text{da+pi}}(\mathbf{x})$  vs.  $\mathcal{H}_{\text{pi}}(\mathbf{x})$ , unless the query point  $\mathbf{x}$  is orthogonal to the variation induced by DA<sup>1</sup>, in which case the size of the interval remains the same. Importantly, Proposition 2 shows that this increase in “sharpness” is a direct consequence of the added variation from of DA in Lemma 1. And because this variation is exogenous and independent of  $\xi$ , our hypothesis/assumption about the strength of confounding  $\Gamma$  in the system should remain the same. This combination of increased data variation, but same confounding assumptions is what reduces ambiguity in PI, resulting in the sharper bounds of Proposition 2. Lastly, in the “large DA” regime, the ellipsoid  $\mathcal{H}_{\text{da+pi}}$  collapses onto  $\text{null}(\Delta)$ , and the interval width  $|\mathcal{H}_{\text{da+pi}}(\mathbf{x})|$  becomes  $\|\Pi_{\Delta}^{\perp} \mathbf{x}\|_{\Sigma_X^{-1}}^2$ . Meaning, the DA removes *all but* the uncertainty that it cannot “see” within its range( $\Delta$ ).

Although smaller identified sets/ intervals are in general desirable, size alone may not be the most appropriate measure of “goodness” of the identified set. The next result is based on worst-case error.

**Theorem 1** (robust bounds with DA). *For  $\mathcal{G}$ -inv.  $\mathbf{f}$ , Assumptions 1 to 3,  $\kappa := \lambda_{\max}(\Sigma_X \Sigma_{GX}^{-1}) \leq 1$ ,*

$$\begin{aligned} E_{\text{worst}}^{\text{do}(X)}(\mathcal{Q}_{\text{pi}}) &= \left( \|\mathbf{h}_{\text{erm}} - \mathbf{f}\|_{\Sigma_X} + r(\Gamma) \right)^2, \\ &\stackrel{(i), (ii)}{\geq} \left( \underbrace{\|\mathbf{h}_{\text{da+erm}} - \mathbf{f}\|_{\Sigma_X}}_{\text{lower estimation error}} + \underbrace{\sqrt{\kappa} \cdot r(\Gamma)}_{\text{sharper bounds}} \right)^2 \stackrel{(ii)}{\geq} E_{\text{worst}}^{\text{do}(X)}(\mathcal{Q}_{\text{da+pi}}). \end{aligned}$$

*Equality iff (i) DA adds low variance  $\kappa = 1$ , and (ii) DA orthogonal to confounding  $\Delta \perp \Sigma_{X,\xi}$ . Also,*

$$\mathbb{E}_{\mathbf{x}} \left[ E_{\text{worst}}^{\text{do}(\mathbf{x})}(\mathcal{Q}_{\text{pi}}) \right] > \underbrace{\|\mathbf{h}_{\text{da+erm}} - \mathbf{f}\|_{\Sigma_X}^2}_{\text{lower estimation error}} + \underbrace{\nu \cdot r(\Gamma)^2}_{\text{sharper bounds}} + s = \mathbb{E}_{\mathbf{x}} \left[ E_{\text{worst}}^{\text{do}(\mathbf{x})}(\mathcal{Q}_{\text{da+pi}}) \right],$$

where  $\nu := \text{tr}(\Sigma_X \Sigma_{GX}^{-1}) < \text{tr}(\Sigma_X \Sigma_X^{-1}) = m$ , queries  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma_X)$  and some slack  $s \geq 0$ .

*Proof.* See Appendix A.1 for the proof.  $\square$

Theorem 1 shows that DA dominates PI on worst-case error through two mechanisms: (i) From Proposition 1, confounding aligned DA causes the PI centroid  $\mathbf{h}_{\text{erm}}$  to drift closer to  $\mathbf{f}$ , bringing the bounds with it, thereby reducing worst-case error. (ii) Independently, from Proposition 2, the bounds themselves shrink, pushing the worst-case point closer still to  $\mathbf{f}$ . Given that the worst-case error bounds how bad the performance of any one hypothesis in the identified set may be, application of a DA pre-processing to PI therefore makes subsequent decision making more robust and reliable. Theorem 1 also gives a lower bound on this improvement via the factor  $\kappa$ , which in the “large DA” regime approaches 0 when  $\Delta$  has full span on  $\mathbb{R}^m$ , but is 1 otherwise. In our linear setting of Assumption 2, the former implies a trivial  $\mathbf{f}$ , and so the last inequality in Theorem 1 more clearly shows the independent, and strictly positive (average) effect  $\nu$  of sharper bounds for individual queries  $\mathbf{x} \sim \mathbb{P}_X$ . Which in the “large DA” regime shrinks to  $\nu \rightarrow \dim(\text{null}(\Delta)) =: k$ , reducing by a factor  $(m - k)/m < 1$ , and directly improving (average) worst-case error for a random query  $\mathbf{x}$ .

#### 4.2.2 VALID BOUNDS WITH DATA AUGMENTATION

Finally, we discuss perhaps the most important property in PI—bound validity. We address this as:

**Theorem 2** (valid bounds with DA). *For any  $\mathcal{G}$ -invariant  $\mathbf{f}$ , it holds under Assumptions 1 to 3 that*

$$E_{\text{approx}}^{\text{do}(X)}(\mathcal{Q}_{\text{da+pi}}) \leq E_{\text{approx}}^{\text{do}(X)}(\mathcal{Q}_{\text{pi}}), \quad \text{equality iff} \quad \mathbb{P} \in \mathcal{Q}_{\text{pi}}, \quad \text{or} \quad \Delta \perp \Sigma_{X,\xi}.$$

*Proof.* See Appendix A.2 for the proof.  $\square$

Meaning the identified set  $\mathcal{H}_{\text{da+pi}}$  is no farther from  $\mathbf{f}$  compared to the original set  $\mathcal{H}_{\text{pi}}$ , and is strictly closer to  $\mathbf{f}$  so long as the DA induced variation aligns with confounding. Of course it follows that when  $\mathcal{Q}_{\text{pi}}$  contains the true joint distribution  $\mathbb{P}$ , then  $\mathbf{f} \in \mathcal{H}_{\text{pi}}$  and so we should also have  $\mathbf{f} \in \mathcal{H}_{\text{da+pi}}$ . Instead of such a simple set inclusion criteria, we keep the more general approximation error framing of Theorem 2 because we also position DA as a tool for improved sensitivity analysis

<sup>1</sup>Intuitively, this would be like rotating an image  $\mathbf{x}$  of a centered circle—the rotation leaves  $\mathbf{x}$  unchanged.

where the constraint set may not necessarily be valid for some values of  $\Gamma$ . Theorem 2 is then reassures that with  $\mathcal{G}$ -invariant  $\mathbf{f}$ , DA at the very least should not cause  $\mathcal{H}_{\text{pi}}$  to drift away from  $\mathbf{f}$ .

Immediately following from Theorem 2, when  $\mathbb{P} \in \mathcal{Q}_{\text{pi}}$ , we also get  $\mathbf{f}^\top \mathbf{x} \in \mathcal{H}_{\text{da+pi}}(\mathbf{x}), \mathcal{H}_{\text{pi}}(\mathbf{x})$ . It is difficult, however, to show a similar result as Theorem 1 for the point-wise evaluation of  $E_{\text{approx}}^{\text{do}(\mathbf{x})}(\mathcal{Q}_{\text{da+pi}})$  vs.  $E_{\text{approx}}^{\text{do}(\mathbf{x})}(\mathcal{Q}_{\text{pi}})$  when  $\mathbb{P} \notin \mathcal{Q}_{\text{pi}}$ , as the approximation error non-trivially depends on the alignment of unknown confounding  $\Sigma_{X,\xi}$  with the query  $\mathbf{x}$ , and both can be arbitrary.

## 5 RELATED WORK

**PI and sensitivity analysis.** We give an account of related PI and sensitivity analysis literature in Section 2.3. Our work is largely orthogonal but complementary to this: we introduce a new source of constraints—symmetry knowledge—that is compatible and composes with existing PI frameworks.

**Symmetry and invariance in causal inference.** Invariance is fundamental to causality: causal mechanisms yield predictions invariant to interventions (Peters et al., 2016). Methods enforce such invariances using auxiliary variables for identification (Peters et al., 2016; Heinze-Deml et al., 2018; Arjovsky et al., 2019; Dance & Bloem-Reddy, 2024; Kilbertus et al., 2020; Singh et al., 2019; Zhang et al., 2023) or robust prediction (Rothenhäusler et al., 2021; Krueger et al., 2021; Christiansen et al., 2022). Akbar et al. (2025) leverage symmetry knowledge for robust prediction; whereas we address the orthogonal, but more principled problem of PI when point identification is infeasible.

**Counterfactual DA.** The literature on counterfactual DA has been the main focus of causal analysis of DA. These methods achieve robust predictors by synthesizing counterfactual examples (Ilse et al., 2021; Yuan et al., 2024; Feder et al., 2023; Pitis et al., 2022; Armengol Urpí et al., 2024; Mahajan et al., 2021; Aloui et al., 2023), but require restrictive assumptions: full SEMs (Yuan et al., 2024; Feder et al., 2023), specific auxiliary variables (Ilse et al., 2021; Feder et al., 2023; Mahajan et al., 2021; Aloui et al., 2023), or complete causal graphs (Pitis et al., 2022; Armengol Urpí et al., 2024). We, like Akbar et al. (2025), require a more accessible symmetry knowledge about the data.

## 6 EXPERIMENTS

We validate our frameworks in finite samples. We fix the augmented sample size to match the original to show that bound sharpening stems from symmetry constraints rather than variance reduction.

### 6.1 SIMULATION EXPERIMENT

We follow Akbar et al. (2025) to instantiate a simulation for the linear Gaussian SEM from Assumption 2. To do this, we first sample the SEM parameters—standard normal matrix  $\mathbf{T} \in \mathbb{R}^{m \times m}$ , and vectors  $\mathbf{f}, \mathbf{e} \in \mathbb{R}^m$ , keeping them fixed throughout the experiment. Then sample standard normal exogenous variables  $(U, N_X, N_Y)$  and pass them through the following model to get observable  $(X, Y)$  confounded by the unobserved  $U$  as:

$$X := \mathbf{T}^\top U + 0.1 \cdot N_X, \quad Y := \mathbf{f}^\top X + \mathbf{e}^\top U + 0.1 \cdot N_Y.$$

Next, we construct a DA  $G$  such that  $\mathbf{f}$  respects  $\mathcal{G}$ -invariance. As with Akbar et al. (2025), we do this by first taking the SVD of  $\mathbf{f}$ ,

$$\mathbf{f} = [\mathbf{u} \quad \mathbf{U}_0] \begin{bmatrix} \sigma & \mathbf{0}_{1 \times (m-1)} \\ \mathbf{0}_{(m-1) \times 1} & \mathbf{0}_{(m-1) \times (m-1)} \end{bmatrix} \begin{bmatrix} \mathbf{v}^\top \\ \mathbf{V}_0^\top \end{bmatrix}.$$

The matrix  $\mathbf{V}_0$  represents the orthonormal basis of  $\text{null}(\mathbf{f})$ . We take  $k$  of these rows to construct  $\mathbf{A} \in \mathbb{R}^{k \times m}$  which defines  $G$ :

$$GX := X + \mathbf{a} \cdot \mathbf{A}^\top G, \quad G \sim \mathcal{N}(\mathbf{0}_k, \mathbf{I}_k).$$

Therefore, by construction we have  $\mathcal{G}$ -invariance and therefore  $\mathbf{f}^\top X = \mathbf{f}^\top GX$ . Figure 2 provides an intuitive visualization

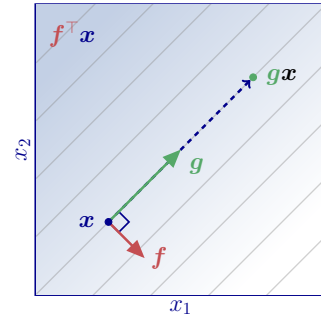


Figure 2: A cartoon visualization of the linear simulation setup. The augmentation generates  $gx$  by randomly translating  $x$  along the level-sets (contours) defined by the causal parameter  $\mathbf{f}$ , using additive noise sampled from the null-space of  $\mathbf{f}$ .



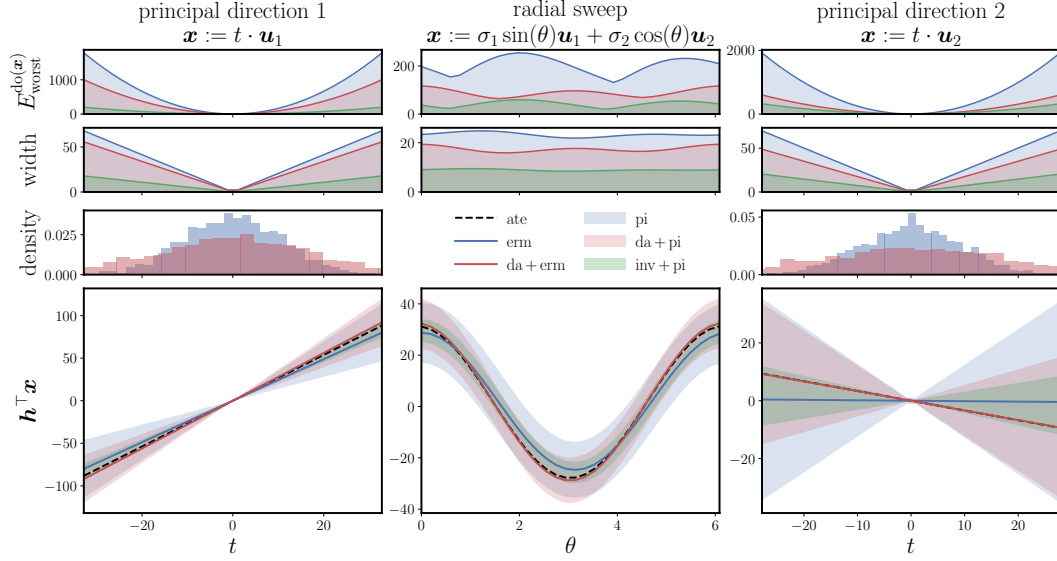


Figure 3: Data augmentation consistently sharpens partial identification bounds in a linear simulation. Across query points aligned with principal components (PC1, PC2) and a radial sweep, DA+PI (red) yields narrower intervals, lower worst-case excess risk ( $E_{\text{worst}}$ ), and predictions closer to the true average treatment effect (ATE, black dash) compared to baseline PI (blue).

of such a transformation. While this construction technically utilizes the ground truth  $\mathbf{f}$ , we treat access to  $\mathbf{A}$  as representing prior domain knowledge about the functional symmetries, noting that this information alone is insufficient to identify  $\mathbf{f}$  due to the unobserved confounding  $\mathbf{e}^\top \mathbf{U}$ .

Taking  $m = 32$ ,  $k = 31$ , we generate  $n = 2048$  samples for  $(X, Y, G)$  with the DA std parameter  $a$  specified as 4. For ERM we use a closed form linear OLS solution. And for PI we use the partial R-squared sensitivity model from Assumption 3 for a range of query points  $\mathbf{x}_0, \mathbf{x}_1$  with the sensitivity parameter set as  $\Gamma = \Gamma_0 = 2^9$ , and  $\epsilon = 2^{-3}$  for the invariance error constraint.

To visualize the results, we chose  $\mathbf{x} := t \cdot \mathbf{u}_1$  and  $\mathbf{x} := t \cdot \mathbf{u}_2$  where  $\mathbf{u}_1, \mathbf{u}_2$  are the first and second principal components of the data. We then sweep  $t$  over  $\pm 3$  standard-deviations, computing intervals  $\mathcal{H}_{\text{pi}}(\mathbf{x}), \mathcal{H}_{\text{da+pi}}(\mathbf{x})$  via convex programming (separately for the upper and lower bounds). The results are shown in Fig. 3 (left, right). Fig. 3 (center) also shows a radial sweep over  $\theta \in [0, 2\pi]$  to generate queries  $\mathbf{x} := \sigma_0 \cdot \sin(\theta) \cdot \mathbf{u}_0 + \sigma_1 \cdot \cos(\theta) \cdot \mathbf{u}_1$ .

## 6.2 OPTICAL DEVICE EXPERIMENT

We utilize the benchmark dataset provided by Janzing & Schölkopf (2018b), consisting of  $3 \times 3$  pixel images  $X$  displayed on a laptop screen which generate voltage readings  $Y$  across a photodiode. The system involves a physically instantiated hidden confounder  $U$  that controls the intensity of two LEDs; the first affects the webcam capturing  $X$ , while the second influences the photodiode measuring  $Y$ . We derive the ground-truth causal predictor  $\mathbf{f}$  by regressing  $Y$  on the joint features  $(\phi(X), U)$ , where  $\phi(X)$  denotes polynomial features of  $X$ . We select the polynomial degree  $d \in \{1, \dots, 5\}$  that best explains the data (degree 2 in our case) and subsequently remove the learned component corresponding to  $U$  to recover  $\mathbf{f}$ . Our choice of DA on  $X$  includes additive Gaussian noise  $G \sim \mathcal{N}(\mathbf{0}, \Sigma_X/10)$ , random vertical/horizontal flips and  $90^\circ$  rotations for DA. We then compute the features  $\phi(GX)$  to be used with PI, setting  $\Gamma = \Gamma_0 = 10^2$  for the partial R-squared model from Assumption 3, and  $\epsilon = 2^{-3}$  for the invariance error constraint on a datasets of  $n = 1000$  samples. Figure 4 shows that DA+PI sharpens bounds over the PI baseline. The visualization approach in the same as in Section 6.1.

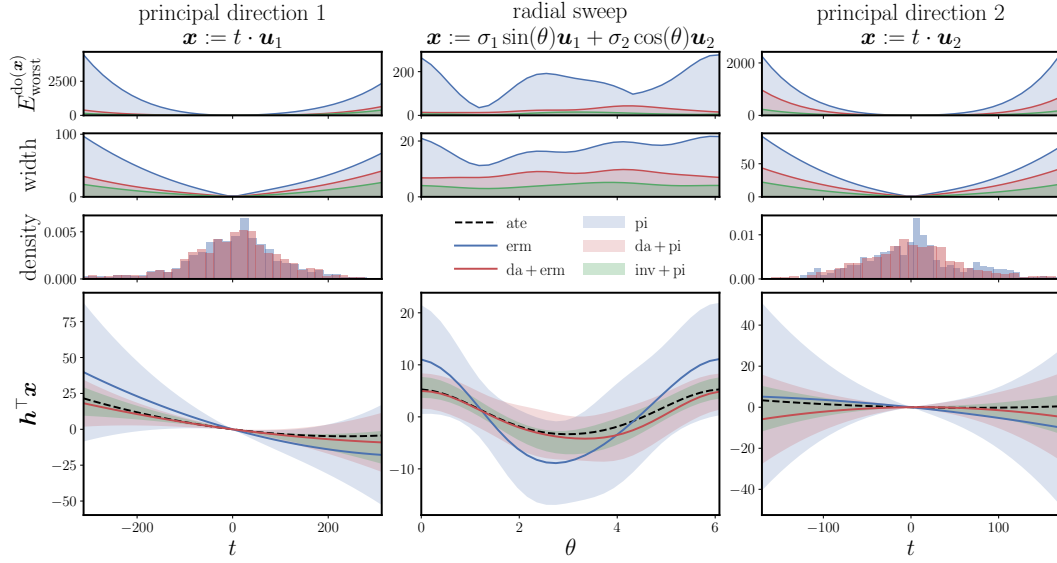


Figure 4: Our method sharpens causal bounds on the real-world Optical Device dataset. Even with complex, non-linear relationships, applying outcome-invariant DA (red) substantially narrows the partial identification bounds compared to the baseline (blue).

## 7 LIMITATIONS, ASSUMPTIONS AND FUTURE WORK

**Symmetry knowledge.** Our approach hinges on the untestable assumption that the target  $f$  is  $\mathcal{G}$ -invariant for the chosen symmetry transformation. While this does require prior knowledge, our framework also allows to handle symmetry miss-specification in the explicit invariance error constraint from Section 4.1 via the  $\epsilon$  parameter. Additionally, we remind the readers that untestable assumptions are fundamental for making any causal conclusions from observational data with unobserved confounding (Pearl, 2009), as is the norm in partial identification. This also includes access to auxiliary variables since the conditional independences that they represent are also merely untestable assumptions. Furthermore, Akbar et al. (2025) argues that a symmetry-based knowledge assumption is actually quite practical given its precedence in the DA and invariance literature (Chen et al., 2020; Lyle et al., 2020; Shao et al., 2022; Fawzi & Frossard, 2015; Dubois et al., 2021; Petrache & Trivedi, 2023; Montasser et al., 2024; Romero & Lohit, 2022; Zhu et al., 2021; Wong et al., 2016).

**Additional covariates.** Many works in PI and sensitivity analysis leverage access to additional auxiliary variables, such as instrumental variables (IVs) and observable confounders or back-doors (Kilbertus et al., 2020; Padh et al., 2023). Even though we do not explicitly model these to keep our analysis simple and tractable, we argue that our symmetry transformation framing is still compatible with them—for example, applying DA on  $X$  does not invalidate an IV that enters into  $X$ .

**Additional partial identification approaches.** Many sensitivity and PI models can be reduced to the constraints in Assumption 3. These include, of course, the partial R-squared model, Rosenbaum (2002), MSM (under a mild bounded marginal ratio assumption), as well as DRO, Wasserstein, total-variation approaches. While a rigorous analysis is left for future work, it is important to specify that our results here are more general than just the partial R-squared model.

## 8 CONCLUSION

We show that causal symmetries sharpen partial identification bounds by restricting the hypothesis space. We operationalize this via explicit invariance constraints and implicit data augmentation. Through construction and linear analysis, respectively, we prove these methods yield valid, strictly tighter, and more robust bounds. Empirically validated and broadly compatible, our framework establishes symmetry as a powerful resource within the tool-belt for causal partial identification.

## ETHICS STATEMENT

The authors have read and adhered to the ICLR Code of Ethics. This work is primarily theoretical and methodological, focusing on the mathematical foundations of using data augmentation for partial identification. The experimental validation relies on a synthetic dataset generated for illustrative purposes and a standard, publicly available benchmark dataset (Optical Device). No human subjects were involved in this research, no new data was collected, and therefore, no Institutional Review Board (IRB) approval was required. The goal of this research is to improve the rigor and reliability of causal inference from observational data, which can lead to more robust and fair decision-making in various applications. We do not foresee any direct negative ethical implications or societal consequences stemming from this work.

## REPRODUCIBILITY STATEMENT

We are committed to ensuring the reproducibility of our work. All theoretical claims made in this paper are supported by detailed, step-by-step proofs, which can be found in the Appendix. The experimental setup for both the simulation study and the real-data experiment is described in Section 6. The complete source code to reproduce all experiments, figures, and results is included as supplementary material with this submission. The code is commented and contains all necessary implementation details, including hyperparameter settings and the specific data generation process for the simulation.

## REFERENCES

- Uzair Akbar, Niki Kilbertus, Hao Shen, Krikamol Muandet, and Bo Dai. An analysis of causal effect estimation using outcome invariant data augmentation. In *Advances in Neural Information Processing Systems*, volume 38. NeurIPS, 2025.
- Ahmed Aloui, Juncheng Dong, Cat P. Le, and Vahid Tarokh. Counterfactual data augmentation with contrastive learning, 2023. arXiv:2311.03630.
- Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization, 2019. arXiv:1907.02893.
- Núria Armengol Urpí, Marco Bagatella, Marin Vlastelica, and Georg Martius. Causal action influence aware counterfactual data augmentation. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235, pp. 1709–1729. PMLR, 2024.
- Alexander Balke and Judea Pearl. Bounds on treatment effects from studies with imperfect compliance. *Journal of the American Statistical Association*, 92(439):1171–1176, 1997.
- Michael M. Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric deep learning: Grids, groups, graphs, geodesics, and gauges, 2021. URL <https://arxiv.org/abs/2104.13478>.
- Gavin Brown and Riccardo Ali. Bias/variance is not the same as approximation/estimation. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=4TnFbvl6hK>.
- Shuxiao Chen, Edgar Dobriban, and Jane H. Lee. A group-theoretic framework for data augmentation. *Journal of Machine Learning Research*, 21(245):1–71, 2020.
- Rune Christiansen, Niklas Pfister, Martin Emil Jakobsen, Nicola Gnecco, and Jonas Peters. A causal framework for distribution generalization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(10):6614–6630, 2022. doi: 10.1109/TPAMI.2021.3094760.
- Carlos Cinelli and Chad Hazlett. Making sense of sensitivity: Extending omitted variable bias. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 82(1):39–67, 12 2019. ISSN 1369-7412. doi: 10.1111/rssb.12348. URL <https://doi.org/10.1111/rssb.12348>.

- Carlos Cinelli and Chad Hazlett. Making sense of sensitivity: Extending omitted variable bias. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 82(1):39–67, 2020.
- Taco Cohen and Max Welling. Group equivariant convolutional networks. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pp. 2990–2999, New York, New York, USA, 20–22 Jun 2016. PMLR.
- Hugh Dance and Benjamin Bloem-Reddy. Causal inference with cocycles, 2024. arXiv:2405.13844.
- Y. Dubois et al. Lossy compression for lossless prediction. In *NeurIPS*, 2021.
- Qiuling Fan, Minghao Li, and Zhengling Wang. Distributionally robust optimization approaches for causal inference with unobserved confounding. *Journal of Machine Learning Research*, 25(87): 1–45, 2024.
- A. Fawzi and P. Frossard. Manitest: Are classifiers really invariant? In *BMVC*, 2015.
- Amir Feder, Yoav Wald, Claudia Shi, Suchi Saria, and David Blei. Data augmentations for improved (large) language model generalization. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 70638–70653. Curran Associates, Inc., 2023.
- Dennis Frauen, Fergus Imrie, Alicia Curth, Valentyn Melnychuk, Stefan Feuerriegel, and Mihaela van der Schaar. A neural framework for generalized causal sensitivity analysis. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=ikX6Dl0Mlc>.
- Florian Gunsilius. A path-sampling method to partially identify causal effects in instrumental variable models, 2020. URL <https://arxiv.org/abs/1910.09502>.
- Wenshuo Guo, Mingzhang Yin, Yixin Wang, and Michael Jordan. Partial identification with noisy covariates: A robust optimization approach. In *First Conference on Causal Learning and Reasoning*, 2022. URL <https://openreview.net/forum?id=-NVBxy0TdU>.
- Christina Heinze-Deml, Jonas Peters, and Nicolai Meinshausen. Invariant causal prediction for nonlinear models. *Journal of Causal Inference*, 6(2), 2018.
- Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, Cambridge, 1985.
- Yaowei Hu, Yongkai Wu, Lu Zhang, and Xintao Wu. A generative adversarial framework for bounding confounded causal effects. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(13):12104–12112, May 2021. doi: 10.1609/aaai.v35i13.17437. URL <https://ojs.aaai.org/index.php/AAAI/article/view/17437>.
- Maximilian Ilse, Jakub M. Tomczak, and Patrick Forré. Selecting data augmentation for simulating interventions. In *International Conference on Machine Learning*, pp. 4555–4562. PMLR, 2021.
- Guido W. Imbens and Charles F. Manski. Confidence intervals for partially identified parameters. *Econometrica*, 72(6):1845–1857, 2004. ISSN 00129682, 14680262. URL <http://www.jstor.org/stable/3598769>.
- Dominik Janzing. Causal regularization. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Dominik Janzing and Bernhard Schölkopf. Detecting non-causal artifacts in multivariate linear regression models. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pp. 2245–2253. PMLR, 2018a.
- Dominik Janzing and Bernhard Schölkopf. Detecting confounding in multivariate linear models via spectral analysis. *Journal of Causal Inference*, 6(1), 2018b.
- Lucas Kania and Ernst Wit. Causal regularization: On the trade-off between in-sample risk and out-of-sample risk guarantees, 2023. arXiv:2205.01593.

- Niki Kilbertus, Matt J. Kusner, and Ricardo Silva. A class of algorithms for general instrumental variable models. In *Advances in Neural Information Processing Systems*, volume 33, pp. 20108–20119, 2020.
- David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghuai Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (REX). In *International Conference on Machine Learning*, pp. 5815–5826. PMLR, 2021.
- Clare Lyle, Mark van der Wilk, Marta Kwiatkowska, Yarin Gal, and Benjamin Bloem-Reddy. On the benefits of invariance in neural networks, 2020. [arXiv:2005.00178](https://arxiv.org/abs/2005.00178).
- Divyat Mahajan, Shruti Tople, and Amit Sharma. Domain generalization using causal matching. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pp. 7313–7324. PMLR, 2021.
- Vahid Balazadeh Meresht, Vasilis Syrgkanis, and Rahul G Krishnan. Partial identification of treatment effects with implicit generative models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=8cUGfg-zUnh>.
- O. Montasser et al. Transformation-invariant learning and theoretical guarantees for ood generalization. In *NeurIPS*, 2024.
- Kirtan Padh, Jakob Zeitler, David Watson, Matt Kusner, Ricardo Silva, and Niki Kilbertus. Stochastic causal programming for bounding treatment effects. In Mihaela van der Schaar, Cheng Zhang, and Dominik Janzing (eds.), *Proceedings of the Second Conference on Causal Learning and Reasoning*, volume 213 of *Proceedings of Machine Learning Research*, pp. 142–176. PMLR, 11–14 Apr 2023. URL <https://proceedings.mlr.press/v213/padh23a.html>.
- Judea Pearl. *Causality*. Cambridge University Press, 2009.
- 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. doi: 10.1111/rssb.12167.
- Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations and learning algorithms*. The MIT Press, 2017.
- M. Petrache and S. Trivedi. Approximation-generalization trade-offs under (approximate) group equivariance. In *NeurIPS*, 2023.
- Silviu Pitis, Elliot Creager, Ajay Mandlekar, and Animesh Garg. MoCoDA: Model-based counterfactual data augmentation. In *Advances in Neural Information Processing Systems*, volume 35, pp. 18143–18156, 2022.
- D. Romero and S. Lohit. Learning partial equivariances from data. In *NeurIPS*, 2022.
- Paul R Rosenbaum. *Observational studies*. Springer, 2002.
- Dominik Rothenhäusler, Nicolai Meinshausen, Peter Bühlmann, and Jonas Peters. Anchor regression: Heterogeneous data meet causality. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 83(2):215–246, 2021.
- Victor Garcia Satorras, Emiel Hoogetboom, and Max Welling. E(n) equivariant graph neural networks, 2022. URL <https://arxiv.org/abs/2102.09844>.
- Jonas Schweisthal, Dennis Frauen, Maresa Schröder, Konstantin Hess, Niki Kilbertus, and Stefan Feuerriegel. Learning representations of instruments for partial identification of treatment effects. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=pgrJPhsk2w>.
- H. Shao et al. A theory of pac learnability under transformation invariances. In *NeurIPS*, 2022.
- Connor Shorten and Taghi M. Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of Big Data*, 6:60, 2019.



- Rahul Singh, Maneesh Sahani, and Arthur Gretton. Kernel instrumental variable regression. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Zhiqiang Tan. A distributional approach for causal inference using propensity scores. *Journal of the American Statistical Association*, 101(476):1619–1637, 2006. doi: 10.1198/016214506000000023. URL <https://doi.org/10.1198/016214506000000023>.
- Chennuru Vankadara, Luca Rendsburg, Ulrike von Luxburg, and Debarghya Ghoshdastidar. Interpolation and regularization for causal learning. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35. Curran Associates, Inc., 2022.
- Vladimir Naumovich Vapnik. *Statistical learning theory*. Wiley, New York, 1998.
- Victor Veitch, Alexander D' Amour, Steve Yadlowsky, and Jacob Eisenstein. Counterfactual invariance to spurious correlations in text classification. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 16196–16208. Curran Associates, Inc., 2021a.
- Victor Veitch, Alexander D' Amour, Steve Yadlowsky, and Jacob Eisenstein. Counterfactual invariance to spurious correlations in text classification. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 16196–16208. Curran Associates, Inc., 2021b.
- S. Wong et al. Understanding data augmentation for classification: When to warp? In *DICTA*, 2016.
- Liyuan Xu and Arthur Gretton. A neural mean embedding approach for back-door and front-door adjustment, 2022. arXiv:2210.06610.
- Fanny Yang, Zuowen Wang, and Christina Heinze-Deml. Invariance-inducing regularization using worst-case transformations suffices to boost accuracy and spatial robustness. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Jianhao Yuan, Francesco Pinto, Adam Davies, and Philip Torr. Not just pretty pictures: Toward interventional data augmentation using text-to-image generators. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=b89JtZj9gm>.
- Rui Zhang, Masaaki Imaizumi, Bernhard Schölkopf, and Krikamol Muandet. Instrumental variable regression via kernel maximum moment loss. *Journal of Causal Inference*, 11(1), 2023.
- S. Zhu et al. Understanding the generalization benefit of model invariance from a data perspective. In *NeurIPS*, 2021.

## A PROOFS

### A.1 PROOF OF THEOREM 1—ROBUST BOUNDS WITH DA

**Theorem 1** (robust bounds with DA). *For  $\mathcal{G}$ -inv.  $\mathbf{f}$ , Assumptions 1 to 3,  $\kappa := \lambda_{\max}(\Sigma_X \Sigma_{GX}^{-1}) \leq 1$ ,*

$$\begin{aligned} E_{\text{worst}}^{\text{do}(X)}(\mathcal{Q}_{\text{pi}}) &= \left( \|\mathbf{h}_{\text{erm}} - \mathbf{f}\|_{\Sigma_X} + r(\Gamma) \right)^2, \\ &\stackrel{(i), (ii)}{\geq} \underbrace{\left( \|\mathbf{h}_{\text{da+erm}} - \mathbf{f}\|_{\Sigma_X} \right)^2}_{\text{lower estimation error}} + \underbrace{\left( \sqrt{\kappa} \cdot r(\Gamma) \right)^2}_{\text{sharper bounds}} \stackrel{(ii)}{\geq} E_{\text{worst}}^{\text{do}(X)}(\mathcal{Q}_{\text{da+pi}}). \end{aligned}$$

*Equality iff (i) DA adds low variance  $\kappa = 1$ , and (ii) DA orthogonal to confounding  $\Delta \perp \Sigma_{X,\xi}$ . Also,*

$$\mathbb{E}_{\mathbf{x}} \left[ E_{\text{worst}}^{\text{do}(\mathbf{x})}(\mathcal{Q}_{\text{pi}}) \right] > \underbrace{\|\mathbf{h}_{\text{da+erm}} - \mathbf{f}\|_{\Sigma_X}^2}_{\text{lower estimation error}} + \underbrace{\nu \cdot r(\Gamma)^2}_{\text{sharper bounds}} + s = \mathbb{E}_{\mathbf{x}} \left[ E_{\text{worst}}^{\text{do}(\mathbf{x})}(\mathcal{Q}_{\text{da+pi}}) \right],$$

where  $\nu := \text{tr}(\Sigma_X \Sigma_{GX}^{-1}) < \text{tr}(\Sigma_X \Sigma_X^{-1}) = m$ , queries  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma_X)$  and some slack  $s \geq 0$ .

*Proof.* We show the two inequalities below in the respective sections.

**Population effect.** Lemma 2, characterizes the identified sets  $\mathcal{H}_{\text{pi}}, \mathcal{H}_{\text{da+pi}}$  as ellipsoids:

$$\mathcal{H}_{\text{pi}} = \left\{ \mathbf{h} \mid \|\mathbf{h} - \mathbf{h}_{\text{erm}}\|_{\Sigma_X}^2 \leq r(\Gamma)^2 \right\}, \quad \mathcal{H}_{\text{da+pi}} = \left\{ \mathbf{h} \mid \|\mathbf{h} - \mathbf{h}_{\text{da+erm}}\|_{\Sigma_{GX}}^2 \leq r(\Gamma)^2 \right\}.$$

Now, from the definition of worst-case excess error in Section 2.3 it follows

$$\begin{aligned} E_{\text{worst}}^{\text{do}(X)}(\mathcal{Q}_{\text{pi}}) &= \max_{\mathbf{Q} \in \mathcal{Q}_{\text{pi}}} E^{\text{do}(X)}(\mathbf{h}_{\text{adj}}^{\mathbf{Q}}), \\ &= \max_{\mathbf{h} \in \mathcal{H}_{\text{pi}}} E^{\text{do}(X)}(\mathbf{h}), \quad (\text{Re-parameterizing in terms of } \mathcal{H}_{\text{pi}}.) \\ &= \max_{\mathbf{h} \in \mathcal{H}_{\text{pi}}} \|\mathbf{h} - \mathbf{f}\|_{\Sigma_X}^2, \\ &= \left( \|\mathbf{h}_{\text{erm}} - \mathbf{f}\|_{\Sigma_X} + r(\Gamma) \right)^2, \quad (\text{Lemma 3}) \end{aligned}$$

where  $r(\Gamma)$  is some constant entirely determined by  $\Gamma$ . Now, we do a similar exercise with  $\mathcal{Q}_{\text{da+pi}}$ ,

$$\begin{aligned} E_{\text{worst}}^{\text{do}(X)}(\mathcal{Q}_{\text{da+pi}}) &= \max_{\mathbf{Q} \in \mathcal{Q}_{\text{da+pi}}} E^{\text{do}(X)}(\mathbf{h}_{\text{adj}}^{\mathbf{Q}}), \\ &= \max_{\mathbf{h} \in \mathcal{H}_{\text{da+pi}}} E^{\text{do}(X)}(\mathbf{h}), \quad (\text{Re-parameterizing in terms of } \mathcal{H}_{\text{da+pi}}.) \\ &= \max_{\mathbf{h} \in \mathcal{H}_{\text{da+pi}}} \|\mathbf{h} - \mathbf{f}\|_{\Sigma_X}^2, \\ &\stackrel{(\star)}{\leq} \left( \|\mathbf{h}_{\text{da+erm}} - \mathbf{f}\|_{\Sigma_X} + r(\Gamma) \cdot \sqrt{\lambda_{\max}(\Sigma_X \Sigma_{GX}^{-1})} \right)^2, \\ &\quad (\text{Lemma 3, = iff } (\mathbf{h}_{\text{da+erm}} - \mathbf{f}) \parallel \mathbf{v}_{\max}(\Sigma_X \Sigma_{GX}^{-1}).) \\ &\stackrel{(\dagger)}{\leq} \left( \|\mathbf{h}_{\text{da+erm}} - \mathbf{f}\|_{\Sigma_X} + r(\Gamma) \right)^2, \quad (\text{Lemma 1, } \lambda_{\max}(\Sigma_X \Sigma_{GX}^{-1}) \leq 1.) \\ &\stackrel{(\ddagger)}{\leq} \left( \|\mathbf{h}_{\text{erm}} - \mathbf{f}\|_{\Sigma_X} + r(\Gamma) \right)^2 = E_{\text{worst}}^{\text{do}(X)}(\mathcal{Q}_{\text{pi}}). \quad (\text{Proposition 1, = iff } \Delta \perp \Sigma_{X,\xi}.) \end{aligned}$$

**Condition for equality.** The bound involves three inequalities:  $(\star)$  the geometric bound from Lemma 3,  $(\dagger)$  the inflated variance implication from Lemma 1, and  $(\ddagger)$  the estimation bound from Proposition 1. Of these,  $(\dagger)$  is immediate, so we investigate  $(\star)$ ,  $(\ddagger)$  in isolation with  $\lambda_{\max} = 1$ . Now, assuming further that the condition  $\Delta \perp \Sigma_{X,\xi}$  is satisfied for  $(\ddagger)$ , it follows from Proposition 1 that

$$\mathbf{h}_{\text{da+erm}} = \mathbf{h}_{\text{erm}} \iff \Delta \perp \Sigma_{X,\xi}.$$

So we go ahead and substitute  $\mathbf{h}_{\text{da+erm}}$  with  $\mathbf{h}_{\text{erm}}$  in  $(\star)$ . From Lemma 3, equality holds iff the bias vector  $(\mathbf{h}_{\text{da+erm}} - \mathbf{f})$ , now  $(\mathbf{h}_{\text{erm}} - \mathbf{f})$ , is a dominant eigenvector of  $\Sigma_X \Sigma_{GX}^{-1}$ . Because  $(\mathbf{h}_{\text{erm}} - \mathbf{f}) = \Sigma_X^{-1} \Sigma_{X,\xi}$  (OLS closed-form), and  $\Sigma_{GX} = \Sigma_X + \Delta$  (Lemma 1), we follow the steps of Propositions 1 and 2 to do a change of basis by jointly diagonalizing  $\Sigma_X, \Delta$  (Lemma 7) to show

$$\Sigma_X \Sigma_{GX}^{-1} (\mathbf{h}_{\text{erm}} - \mathbf{f}) = \Sigma_X (\Sigma_X + \Delta)^{-1} (\Sigma_X^{-1} \Sigma_{X,\xi}) = (\mathbf{h}_{\text{erm}} - \mathbf{f}), \quad \Longleftrightarrow \quad \Delta \perp \Sigma_{X,\xi}.$$

The bias  $(\mathbf{h}_{\text{erm}} - \mathbf{f})$  is an eigenvector of  $\Sigma_X \Sigma_{GX}^{-1}$  with eigenvalue  $\lambda_{\max} = 1$ . Therefore, equality holds for both  $(\star)$  and  $(\dagger)$  iff  $\Delta \perp \Sigma_{X,\xi}$  and the residual improvement is solely from radius contraction  $\lambda_{\max} < 1$ . Conditions for  $(\star)$ ,  $(\dagger)$ ,  $(\ddagger)$  form conditions (i), (ii) in the statement.

**(Average) individual effect.** Define  $J(\mathcal{Q}) := \mathbb{E}_{\mathbf{x}} \left[ E_{\text{worst}}^{\text{do}(\mathbf{x})}(\mathcal{Q}) \right]$  for queries  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma_X)$ . From Lemma 2, the worst-case risk at a query point  $\mathbf{x}$  is (bias + radius) squared:

$$E_{\text{worst}}^{\text{do}(\mathbf{x})}(\mathcal{Q}_{\text{pi}}) = \left( |(\mathbf{h}_{\text{erm}} - \mathbf{f})^\top \mathbf{x}| + r(\Gamma) \|\mathbf{x}\|_{\Sigma_X^{-1}} \right)^2.$$

Expanding the square, we decompose the total expected risk for DA into three terms:

$$\begin{aligned} J(\mathcal{Q}_{\text{da+pi}}) &= \underbrace{\mathbb{E}_{\mathbf{x}} \left[ \|\mathbf{h}_{\text{da+erm}} - \mathbf{f}\|_{\mathbf{x}}^2 \right]}_{\text{(a) estimation error}} + \underbrace{r(\Gamma)^2 \mathbb{E}_{\mathbf{x}} \left[ \|\mathbf{x}\|_{\Sigma_{GX}^{-1}}^2 \right]}_{\text{(b) average radius}} \\ &\quad + \underbrace{2r(\Gamma) \mathbb{E}_{\mathbf{x}} \left[ |(\mathbf{h}_{\text{da+erm}} - \mathbf{f})^\top \mathbf{x}| \|\mathbf{x}\|_{\Sigma_{GX}^{-1}} \right]}_{\text{(c) interaction of (a), (b)}}. \end{aligned}$$

We analyze the reduction  $J(\mathcal{Q}_{\text{pi}}) - J(\mathcal{Q}_{\text{da+pi}})$  term by term:

- a) **Estimation error:** DA+ERM dominates ERM from Proposition 1, equality iff  $\Delta \perp \Sigma_{X,\xi}$
- b) **Average radius:** DA+PI *strictly* dominates PI from Proposition 2, as  $\mathbf{x} \not\perp \Delta$  almost surely. Also, expand  $\|\mathbf{x}\|_{\Sigma_{GX}^{-1}}^2$  into a trace term, and then use its cyclic permutation invariance gets

$$\mathbb{E}_{\mathbf{x}} \left[ \|\mathbf{x}\|_{\Sigma_{GX}^{-1}}^2 \right] = \text{tr}(\Sigma_{GX}^{-1} \Sigma_X) < \text{tr}(\Sigma_X^{-1} \Sigma_X) = m.$$

- c) **Interaction term:** *Strictly* lower for any non-trivial DA  $\Delta \neq \mathbf{0}$  from Lemma 5.

Concluding that:

$$J(\mathcal{Q}_{\text{da+pi}}) < J(\mathcal{Q}_{\text{pi}}) \quad \text{when} \quad \Delta \neq \mathbf{0}.$$

□

## A.2 PROOF OF THEOREM 2—VALID BOUNDS WITH DA

**Theorem 2** (valid bounds with DA). *For any  $\mathcal{G}$ -invariant  $\mathbf{f}$ , it holds under Assumptions 1 to 3 that*

$$E_{\text{approx}}^{\text{do}(X)}(\mathcal{Q}_{\text{da+pi}}) \leq E_{\text{approx}}^{\text{do}(X)}(\mathcal{Q}_{\text{pi}}), \quad \text{equality iff } \mathbb{P} \in \mathcal{Q}_{\text{pi}}, \quad \text{or} \quad \Delta \perp \Sigma_{X,\xi}.$$

*Proof.* From Lemma 2, we can characterize the identified sets  $\mathcal{H}_{\text{pi}}$ ,  $\mathcal{H}_{\text{da+pi}}$  as ellipsoids of the form

$$\mathcal{H}_{\text{pi}} = \left\{ \mathbf{h} \mid \|\mathbf{h} - \mathbf{h}_{\text{erm}}\|_{\Sigma_X}^2 \leq r(\Gamma)^2 \right\}, \quad \mathcal{H}_{\text{da+pi}} = \left\{ \mathbf{h} \mid \|\mathbf{h} - \mathbf{h}_{\text{da+erm}}\|_{\Sigma_{GX}}^2 \leq r(\Gamma)^2 \right\}.$$

First consider  $\mathbb{P} \notin \mathcal{Q}_{\text{pi}}$ . Now, from the definition of approximation error in Section 2.3 it follows

$$\begin{aligned} E_{\text{approx}}^{\text{do}(X)}(\mathcal{Q}_{\text{pi}}) &= \min_{\mathbb{Q} \in \mathcal{Q}_{\text{pi}}} E^{\text{do}(X)}(h_{\text{adj}}^{\mathbb{Q}}), \\ &= \min_{\mathbf{h} \in \mathcal{H}_{\text{pi}}} E^{\text{do}(X)}(\mathbf{h}), \quad (\text{Re-parameterizing in terms of } \mathcal{H}_{\text{pi}}.) \\ &= \min_{\mathbf{h} \in \mathcal{H}_{\text{pi}}} \|\mathbf{h} - \mathbf{f}\|_{\Sigma_X}^2, \\ &= \left( \|\mathbf{h}_{\text{erm}} - \mathbf{f}\|_{\Sigma_X} - r(\Gamma) \right)^2, \quad (\text{Lemma 4}) \end{aligned}$$

where  $r(\Gamma)$  is some constant entirely determined by  $\Gamma$ . Now, we do a similar exercise with  $\mathcal{Q}_{\text{da+pi}}$ ,

$$\begin{aligned} E_{\text{approx}}^{\text{do}(X)}(\mathcal{Q}_{\text{da+pi}}) &= \min_{\mathbb{Q} \in \mathcal{Q}_{\text{da+pi}}} E^{\text{do}(X)}(h_{\text{adj}}^{\mathbb{Q}}), \\ &= \min_{\mathbf{h} \in \mathcal{H}_{\text{da+pi}}} E^{\text{do}(X)}(\mathbf{h}), \quad (\text{Re-parameterizing in terms of } \mathcal{H}_{\text{da+pi}}.) \\ &= \min_{\mathbf{h} \in \mathcal{H}_{\text{da+pi}}} \|\mathbf{h} - \mathbf{f}\|_{\Sigma_X}^2, \\ &\stackrel{(\heartsuit)}{\leq} \left( 1 - \frac{r(\Gamma)}{\|\mathbf{h}_{\text{da+erm}} - \mathbf{f}\|_{\Sigma_{GX}}} \right)^2 \|\mathbf{h}_{\text{da+erm}} - \mathbf{f}\|_{\Sigma_X}^2, \\ &\quad (\text{Lemma 4, = iff } (\mathbf{h}_{\text{da+erm}} - \mathbf{f}) \parallel \mathbf{v}(\Sigma_X \Sigma_{GX}^{-1}).) \\ &= \left( \|\mathbf{h}_{\text{da+erm}} - \mathbf{f}\|_{\Sigma_{GX}} - r(\Gamma) \right)^2 \frac{\|\mathbf{h}_{\text{da+erm}} - \mathbf{f}\|_{\Sigma_X}^2}{\|\mathbf{h}_{\text{da+erm}} - \mathbf{f}\|_{\Sigma_{GX}}^2}, \\ &\stackrel{(\spadesuit)}{\leq} \left( \|\mathbf{h}_{\text{erm}} - \mathbf{f}\|_{\Sigma_X} - r(\Gamma) \right)^2 = E_{\text{approx}}^{\text{do}(X)}(\mathcal{Q}_{\text{pi}}), \quad (\text{Similar to Proposition 1, = iff } \Delta \perp \Sigma_{X,\xi}.) \end{aligned}$$

where the last inequality ( $\dagger$ ) follows from a similar approach as used in Proposition 1 to show that

$$\|\mathbf{h}_{\text{da+erm}} - \mathbf{f}\|_{\Sigma_X}^2 \leq \|\mathbf{h}_{\text{da+erm}} - \mathbf{f}\|_{\Sigma_{GX}}^2 \leq \|\mathbf{h}_{\text{erm}} - \mathbf{f}\|_{\Sigma_X}^2,$$

which holds with equality if and only if  $\Delta \perp \Sigma_{X,\xi}$ . The case for  $\mathbb{P} \in \mathcal{Q}_{\text{pi}}$  is trivial from Lemma 4.

**Condition for equality.** The two types of inequalities that comprise the given approximation error bound are Proposition 1-type estimation bias related ( $\spadesuit$ ), and the ellipsoidal geometry inequality ( $\heartsuit$ ) from Lemma 4. We can proceed similar to the corresponding section in Theorem 1 to show that:

$$\Delta \perp \Sigma_{X,\xi} \iff (\mathbf{h}_{\text{da+erm}} - \mathbf{f}) \parallel \mathbf{v}_{\max}(\Sigma_X \Sigma_{GX}^{-1}).$$

Since condition of ( $\heartsuit$ ) requires alignment of  $(\mathbf{h}_{\text{da+erm}} - \mathbf{f})$  with *any* eigenvector  $\mathbf{v}(\Sigma_X \Sigma_{GX}^{-1})$ , therefore  $\mathbf{v}_{\max}$  suffices. Consequently, equality holds for both ( $\heartsuit$ ) and ( $\spadesuit$ ) iff  $\Delta \perp \Sigma_{X,\xi}$ .  $\square$

## A.3 PROOF OF PROPOSITION 1—ESTIMATION WITH DA (AKBAR ET AL. (2025) LIFTED)

**Proposition 1** (estimation with DA (Akbar et al. (2025) lifted)). *For  $\mathcal{G}$ -inv.  $\mathbf{f}$ , Assumptions 1 and 2,*

$$0 \leq \frac{\kappa}{1+\kappa} \cdot \underbrace{\|\Pi_{\Delta}(\mathbf{h}_{\text{erm}} - \mathbf{f})\|_{\Sigma_X}^2}_{\text{estimation error within range}(\Delta)} \leq E^{\text{do}(X)}(\mathbf{h}_{\text{erm}}) - E^{\text{do}(X)}(\mathbf{h}_{\text{da+erm}}),$$

$$\leq \|\Pi_{\Delta}(\mathbf{h}_{\text{erm}} - \mathbf{f})\|_{\Sigma_X}^2, \quad \text{eq. iff} \quad \underbrace{\Delta \perp \Sigma_{X,\xi}}_{\text{DA orthogonal to confounding}},$$

where  $\kappa := \lambda_{\min}^+(\Sigma_X^{-1}\Delta) < \infty$  represents the lowest positive eigenvalue of the product  $\Sigma_X^{-1}\Delta$ .

*Proof.* We start by first investigating the post-DA confounding vector  $\mathbb{E}[(GX)\xi^\top] = \Sigma_{GX,\xi}$  as

$$\begin{aligned} \Sigma_{GX,\xi} &= \mathbb{E}[(GX)\xi^\top] \\ &= \mathbb{E}_{X,\xi}[\mathbb{E}_G[GX | X, \xi]\xi^\top] && \text{(Law of total expectation.)} \\ &= \mathbb{E}_{X,\xi}[\mathbb{E}_G[GX | X]\xi^\top] && (G \text{ exogenous} \implies G \perp\!\!\!\perp \xi | X.) \\ &= \mathbb{E}_{X,\xi}[X\xi^\top] = \Sigma_{X,\xi}. && (\text{As } \mathbb{E}[GX | X] = X \text{ from Assumption 1.}) \end{aligned}$$

Now define  $\mathbf{c} := \Sigma_{X,\xi} = \Sigma_{GX,\xi}$  for brevity. The estimation error in Eq. (6) for the baseline ERM and DA+ERM is governed by the projection of confounding  $\mathbf{c}$  onto the respective data manifolds as:

$$\begin{aligned} E^{\text{do}(X)}(\mathbf{h}_{\text{erm}}) &= \|\Sigma_X^{-1}\mathbf{c}\|_{\Sigma_X}^2 = \mathbf{c}^\top \Sigma_X^{-1}\mathbf{c}, \\ E^{\text{do}(X)}(\mathbf{h}_{\text{da+erm}}) &= \|\Sigma_{GX}^{-1}\mathbf{c}\|_{\Sigma_X}^2 = \mathbf{c}^\top \Sigma_{GX}^{-1}\Sigma_X \Sigma_{GX}^{-1}\mathbf{c} \end{aligned}$$

Using  $\Sigma_X = \Sigma_{GX} - \Delta$  and the Resolvent Identity  $\Sigma_X^{-1} - \Sigma_{GX}^{-1} = \Sigma_{GX}^{-1}\Delta\Sigma_X^{-1}$ , we get:

$$\begin{aligned} E^{\text{do}(X)}(\mathbf{h}_{\text{da+erm}}) &= \mathbf{c}^\top \Sigma_{GX}^{-1}(\Sigma_{GX} - \Delta)\Sigma_{GX}^{-1}\mathbf{c} \\ &= \mathbf{c}^\top \Sigma_{GX}^{-1}\mathbf{c} - \mathbf{c}^\top (\Sigma_{GX}^{-1}\Delta\Sigma_{GX}^{-1})\mathbf{c} \\ &= (\mathbf{c}^\top \Sigma_X^{-1}\mathbf{c} - \mathbf{c}^\top \Sigma_{GX}^{-1}\Delta\Sigma_X^{-1}\mathbf{c}) - \mathbf{c}^\top (\Sigma_{GX}^{-1}\Delta\Sigma_{GX}^{-1})\mathbf{c} \\ &= E^{\text{do}(X)}(\mathbf{h}_{\text{erm}}) - \mathbf{c}^\top (\Sigma_{GX}^{-1}\Delta\Sigma_X^{-1})\mathbf{c} - \mathbf{c}^\top (\Sigma_{GX}^{-1}\Delta\Sigma_{GX}^{-1})\mathbf{c}, \\ &= E^{\text{do}(X)}(\mathbf{h}_{\text{erm}}) \quad \underbrace{0 \leq \text{first-order reduction}}_{\substack{-\mathbf{c}^\top \Sigma_X^{-1}(\Sigma_X \Sigma_{GX}^{-1}\Delta)\Sigma_X^{-1}\mathbf{c} \\ -\mathbf{c}^\top \Sigma_X^{-1}(\Sigma_X \Sigma_{GX}^{-1}\Delta\Sigma_{GX}^{-1}\Sigma_X)\Sigma_X^{-1}\mathbf{c}. \\ 0 \leq \text{second-order reduction}}} \end{aligned}$$

Both reduction terms are quadratic forms of the PSD matrix  $\Delta$  and are therefore non-negative.

Define  $\delta$  as their sum. Lemma 6 lower-bounds the first-order term, and by extension lower-bounds  $\delta$ :

$$0 \leq \frac{\kappa}{1+\kappa} \cdot \|\Pi_{\Delta}(\mathbf{h}_{\text{erm}} - \mathbf{f})\|_{\Sigma_X}^2 \leq \text{first order term} \stackrel{(\blacktriangle)}{\leq} \delta.$$

Trace the same steps as Lemma 6 to bound  $\delta$  from above via the simultaneous basis from Lemma 7 ( $\Sigma_X = S^\top S$ ,  $\Delta = S^\top DS$ ). Taking  $\mathbf{z} := S\Sigma_X^{-1}\mathbf{c}$  and eigenvalues  $D_{ii}$  of  $\Sigma_X^{-1}\Delta$ , we can show

$$\delta = \sum_i z_i^2 \cdot \left( \underbrace{\frac{D_{ii}}{1+D_{ii}}}_{\text{1st order}} + \underbrace{\frac{D_{ii}}{(1+D_{ii})^2}}_{\text{2nd order}} \right) \stackrel{(\nabla)}{\leq} \sum_{i:D_{ii}>0} 1 \cdot z_i^2 = \|\Pi_{\Delta}(\mathbf{h}_{\text{erm}} - \mathbf{f})\|_{\Sigma_X}^2.$$

**Condition for equality.** Equality holds for  $(\blacktriangle)$  iff  $\Delta \perp \Sigma_{X,\xi}$ , as otherwise the second-order term is strictly positive. Equality also holds for  $(\nabla)$  iff  $\Delta \perp \Sigma_{X,\xi}$ , because that entails  $z_i = 0$  whenever  $D_{ii} > 0$  so that the sums on both sides go to 0.  $\square$



## A.4 PROOF OF PROPOSITION 2—SHARPER BOUNDS WITH DA

**Proposition 2** (sharper bounds with DA). *For Assumptions 1 to 3, Lebesgue measure (volume)  $|\cdot|$ ,*

$$\frac{|\mathcal{H}_{\text{da+pi}}|}{|\mathcal{H}_{\text{pi}}|} = \sqrt{\frac{\det \Sigma_X}{\det \Sigma_{GX}}} < 1, \quad \frac{|\mathcal{H}_{\text{da+pi}}(\mathbf{x})|}{|\mathcal{H}_{\text{pi}}(\mathbf{x})|} = \frac{\|\mathbf{x}\|_{\Sigma_{GX}^{-1}}}{\|\mathbf{x}\|_{\Sigma_X^{-1}}} \leq 1, \quad \text{equality iff } \mathbf{x} \perp \Delta.$$

*Proof.* We compare the geometric properties of the identified sets as characterized by Lemma 2.

**Ellipsoid volume (global contraction).** Given that the volume of a  $\Sigma$ -ellipsoid  $\propto (\det \Sigma)^{-1/2}$ , it immediately follows from Lemmas 1 and 2, and the monotonicity of determinant for SPD matrices:

$$\begin{aligned} \Sigma_X \preceq \Sigma_{GX} &\implies \det(\Sigma_X) < \det(\Sigma_{GX}), \\ &\implies \det(\Sigma_{GX})^{-1/2} < \det(\Sigma_X)^{-1/2} \implies |\mathcal{H}_{\text{da+pi}}| < |\mathcal{H}_{\text{pi}}|. \end{aligned}$$

**Interval width (point-wise contraction).** From Lemma 2, the width of the interval  $\mathcal{H}_{\text{pi}}(\mathbf{x})$  is simply  $2r(\Gamma) \cdot \|\mathbf{x}\|_{\Sigma_X^{-1}}$ . It then immediately follows from Lemma 1 and definition of the PSD order:

$$\begin{aligned} \Sigma_X \preceq \Sigma_{GX} &\implies \Sigma_{GX}^{-1} \preceq \Sigma_X^{-1}, \\ &\implies \mathbf{x}^\top \Sigma_{GX}^{-1} \mathbf{x} \leq \mathbf{x}^\top \Sigma_X^{-1} \mathbf{x} \implies |\mathcal{H}_{\text{da+pi}}(\mathbf{x})| \leq |\mathcal{H}_{\text{pi}}(\mathbf{x})|. \end{aligned}$$

**Condition for equality.** The interval width is strictly smaller for  $\mathcal{H}_{\text{da+pi}}(\mathbf{x})$  compared to  $\mathcal{H}_{\text{pi}}(\mathbf{x})$  unless the query point  $\mathbf{x}$  lies in the null space of the difference  $\Delta := \Sigma_{GX} - \Sigma_X$ . From Lemma 7,

$$\Sigma_{GX}^{-1} = (\Sigma_X + \Delta)^{-1} = (S^\top S + S^\top D S)^{-1} = S^{-1}(I + D)^{-1} S^{-\top}.$$

When we analyze the ratio of squared norms using the basis  $\mathbf{z} := S^{-\top} \mathbf{x}$ , it simplifies to:

$$\frac{\|\mathbf{x}\|_{\Sigma_{GX}^{-1}}^2}{\|\mathbf{x}\|_{\Sigma_X^{-1}}^2} = \frac{\mathbf{z}^\top (I + D)^{-1} \mathbf{z}}{\mathbf{z}^\top \mathbf{z}} = \frac{\sum_i z_i^2 (1 + D_{ii})^{-1}}{\sum_i z_i^2}.$$

Since  $D$  is non-negative, the term  $(1 + D_{ii})^{-1} < 1$  whenever  $D_{ii} > 0$ . Therefore, the ratio is strictly less than 1 unless  $\mathbf{z}$  is supported only on indices where  $D_{ii} = 0$ . This requires  $\mathbf{z}^\top D \mathbf{z} = 0$ , which transforms back to the condition that  $\mathbf{x}$  must lie in the null-space of  $\Delta$  (i.e.,  $\mathbf{x} \perp \Delta$ ).  $\square$

## A.5 MISCELLANEOUS SUPPORTING LEMMAS

**Lemma 1** (added exogenous variation with DA). *Under Assumption 1,  $G$  inflates the data variance,*

$$\Delta := \Sigma_{GX} - \Sigma_X \succcurlyeq \mathbf{0}, \quad \text{equality iff } GX = X \text{ a.s.}$$

*Proof.* Represent  $Z := GX$ . Now, by applying the Law of Total Covariance conditioning on  $X$ ,

$$\Sigma_{GX} = \mathbb{E}[\text{Cov}(GX | X)] + \text{Cov}(\mathbb{E}[GX | X]). \quad (7)$$

By Assumption 1 (unbiased group action) we have  $\mathbb{E}[GX | X] = X$ , and the second term reduces to

$$\text{Cov}(\mathbb{E}[GX | X]) = \text{Cov}(X) = \Sigma_X. \quad (8)$$

The first term represents the exogenous variation injected by the group action. Let  $\Delta = \mathbb{E}[\text{Cov}(GX | X)]$ . Since covariance matrices are PSD by definition, we have  $\Delta \succcurlyeq \mathbf{0}$ .

**Condition for equality.** The inequality holds with equality ( $\Sigma_{GX} = \Sigma_X$ ) iff the injected noise matrix  $\Delta = \mathbf{0}$ . Since  $\text{Cov}(GX | X) \succcurlyeq \mathbf{0}$  almost surely, its expectation is zero if and only if  $\text{Cov}(GX | X) = \mathbf{0}$  almost surely. This implies  $GX$  is a deterministic function of  $X$ . Given the unbiased assumption  $\mathbb{E}[GX | X] = X$ , this forces  $GX = X$  almost surely (i.e.,  $G$  acts as identity over support of  $X$ ). Therefore, for any non-trivial augmentation, the inequality  $\Delta \succcurlyeq \mathbf{0}$  is strict.  $\square$

**Lemma 2** (characterizing the identified set in a linear, Gaussian case). *Under Assumptions 2 and 3,*

$$\mathcal{H}_{\text{pi}} = \left\{ \mathbf{h} \mid \|\mathbf{h} - \mathbf{h}_{\text{erm}}\|_{\Sigma_X}^2 \leq r(\Gamma)^2 \right\},$$

where the ellipsoid radius  $r(\Gamma) \geq 0$  depends on the choice of constraint parameters. Furthermore,

$$\mathcal{H}_{\text{pi}}(\mathbf{x}) = \left[ \mathbf{h}_{\text{erm}}^\top \mathbf{x} - r(\Gamma) \cdot \|\mathbf{x}\|_{\Sigma_X^{-1}}, \quad \mathbf{h}_{\text{erm}}^\top \mathbf{x} + r(\Gamma) \cdot \|\mathbf{x}\|_{\Sigma_X^{-1}} \right].$$

*Proof.* Compute the population covariance

$$\text{Cov}(X, Y) = \text{Cov}(X, \mathbf{f}^\top X + \xi) = \Sigma_X \mathbf{f} + \Sigma_{X, \xi},$$

so the (naïve) ERM estimand satisfies

$$\mathbf{h}_{\text{erm}} = \Sigma_X^{-1} \text{Cov}(X, Y) = \mathbf{f} + \Sigma_X^{-1} \Sigma_{X, \xi}.$$

Let  $\mathbf{b} := \mathbf{h}_{\text{erm}} - \mathbf{f} = \Sigma_X^{-1} \Sigma_{X, \xi}$ . By the partial- $R^2$  constraint in Assumption 3

$$R_{\xi|X}^2 = \frac{\Sigma_{X, \xi}^\top \Sigma_X^{-1} \Sigma_{X, \xi}}{\sigma_\xi^2} \leq \Gamma,$$

we have

$$\Sigma_{X, \xi}^\top \Sigma_X^{-1} \Sigma_{X, \xi} \leq \sigma_\xi^2 \Gamma.$$

Substituting  $\Sigma_{X, \xi} = \Sigma_{XX} \mathbf{b} = \Sigma_{XX}(\mathbf{h}_{\text{erm}} - \mathbf{f})$  yields

$$(\mathbf{h}_{\text{erm}} - \mathbf{f})^\top \Sigma_{XX} (\mathbf{h}_{\text{erm}} - \mathbf{f}) \leq \sigma_\xi^2 \Gamma,$$

which is equivalent to

$$\|\mathbf{f} - \mathbf{h}_{\text{erm}}\|_{\Sigma_{XX}}^2 \leq \sigma_\xi^2 \Gamma \leq \Gamma_0 \Gamma.$$

Thus the identified set for  $\mathbf{f}$  is the stated ellipsoid with radius  $r(\Gamma)^2 = \Gamma_0 \Gamma$ . The centred Gaussian assumption guarantees the linear projection interpretation used above is exact.

Lastly, since the identified set is an ellipsoid, maximizing/minimizing a linear functional  $\mathbf{f}^\top \mathbf{x}$  is just moving along its principal axis in the direction of  $\mathbf{x}$ , giving us the bounds for  $\mathcal{H}_{\text{pi}}(\mathbf{x})$ .  $\square$

**Lemma 3** (upper bound on distance of a point to farthest point on ellipsoid). *Take ellipsoid  $\mathcal{O} \subset \mathbb{R}^n$*

$$\mathcal{O} = \left\{ \mathbf{x} \mid (\mathbf{x} - \mathbf{x}_0)^\top \Sigma_0 (\mathbf{x} - \mathbf{x}_0) \leq r_0^2 \right\},$$

*with radius  $r_0$ , centered at  $\mathbf{x}_0$  and shape defined by the SPD matrix  $\Sigma_0 \succ 0$ . For some arbitrary point  $\mathbf{y} \in \mathbb{R}^n$ , denote its distance from the farthest point on  $\mathcal{O}$  as weighted by an SPD  $\Sigma \succ 0$  with*

$$D_{\Sigma}^{\max}(\mathbf{y}, \mathcal{O}) := \max_{\mathbf{x} \in \mathcal{O}} \|\mathbf{y} - \mathbf{x}\|_{\Sigma}.$$

*This distance is upper bounded as follows, with  $\mathbf{v}_{\max}$  as the eigenvector corresponding to  $\lambda_{\max}$ .*

$$D_{\Sigma}^{\max}(\mathbf{y}, \mathcal{O}) \leq \|\mathbf{y} - \mathbf{x}_0\|_{\Sigma} + r_0 \cdot \sqrt{\lambda_{\max}(\Sigma \Sigma_0^{-1})},$$

*equality iff  $\mathbf{y} - \mathbf{x}_0 \parallel \mathbf{v}_{\max}(\Sigma \Sigma_0^{-1})$ .*

*Proof.* By triangle inequality

$$\|\mathbf{y} - \mathbf{x}\|_{\Sigma} \leq \|\mathbf{y} - \mathbf{x}_0\|_{\Sigma} + \|\mathbf{x}_0 - \mathbf{x}\|_{\Sigma}.$$

Now, simply maximizing both sides over  $\mathbf{x} \in \mathcal{O}$ ,

$$\max_{\mathbf{x} \in \mathcal{O}} \|\mathbf{y} - \mathbf{x}\|_{\Sigma} \leq \max_{\mathbf{x} \in \mathcal{O}} (\|\mathbf{y} - \mathbf{x}_0\|_{\Sigma} + \|\mathbf{x}_0 - \mathbf{x}\|_{\Sigma}) = \|\mathbf{y} - \mathbf{x}_0\|_{\Sigma} + \max_{\mathbf{x} \in \mathcal{O}} \|\mathbf{x}_0 - \mathbf{x}\|_{\Sigma}.$$

The last term  $\max_{\mathbf{x} \in \mathcal{O}} \|\mathbf{x}_0 - \mathbf{x}\|_{\Sigma}$  is simply the radius of the ellipsoid in the  $\Sigma$ -norm, which is equal to  $r_0 \cdot \sqrt{\lambda_{\max}(\Sigma \Sigma_0^{-1})}$ . The result follows.

**Condition for equality.** The triangle inequality holds with equality iff  $(\mathbf{y} - \mathbf{x}_0)$  and  $(\mathbf{x} - \mathbf{x}_0)$  are collinear. The second term is maximized when  $(\mathbf{x} - \mathbf{x}_0)$  aligns with the dominant eigenvector  $\mathbf{v}_{\max}(\Sigma \Sigma_0^{-1})$  (the generalized principal axis). Therefore, the total bound is tight iff  $(\mathbf{y} - \mathbf{x}_0)$  is itself an eigenvector corresponding to  $\lambda_{\max}(\Sigma \Sigma_0^{-1})$ , i.e.  $(\mathbf{y} - \mathbf{x}_0) \parallel \mathbf{v}_{\max}(\Sigma \Sigma_0^{-1})$ .  $\square$

**Lemma 4** (upper bound on distance of a point to an ellipsoid). *Take the following ellipsoid  $\mathcal{O} \subset \mathbb{R}^n$*

$$\mathcal{O} = \left\{ \mathbf{x} \mid (\mathbf{x} - \mathbf{x}_0)^\top \Sigma_0 (\mathbf{x} - \mathbf{x}_0) \leq r_0^2 \right\},$$

*with radius  $r_0$ , centered at  $\mathbf{x}_0$  and shape defined by the SPD matrix  $\Sigma_0 \succ 0$ . For some arbitrary point  $\mathbf{y} \in \mathbb{R}^n$ , denote its distance from  $\mathcal{O}$  as weighted by an SPD  $\Sigma \succ 0$  with the following notation*

$$D_{\Sigma}^{\min}(\mathbf{y}, \mathcal{O}) := \min_{\mathbf{x} \in \mathcal{O}} \|\mathbf{y} - \mathbf{x}\|_{\Sigma}.$$

*This distance is upper bounded by the following closed-form, with  $\mathbf{v}(\cdot)$  as any arbitrary eigenvector.*

$$D_{\Sigma}^{\min}(\mathbf{y}, \mathcal{O}) \leq \begin{cases} 0, & \mathbf{y} \in \mathcal{O}, \\ \left(1 - \frac{r_0}{\|\mathbf{y} - \mathbf{x}_0\|_{\Sigma_0}}\right) \|\mathbf{y} - \mathbf{x}_0\|_{\Sigma}, & \mathbf{y} \notin \mathcal{O}, \end{cases}$$

*equality iff  $\mathbf{y} \in \mathcal{O}$ , or  $\mathbf{y} - \mathbf{x}_0 \parallel \mathbf{v}(\Sigma \Sigma_0^{-1})$ .*

*Proof.* The result for  $\mathbf{y} \in \mathcal{O}$  case is immediate. To show the bound for  $\mathbf{y} \notin \mathcal{O}$ , consider the ray

$$\mathbf{x}(r) := \mathbf{x}_0 + r \cdot (\mathbf{y} - \mathbf{x}_0), \quad r \in [0, 1],$$

going from the ellipsoid center  $\mathbf{x}_0$  through  $\mathbf{y}$ . This ray intersects with the ellipsoid boundary at

$$r^* = \frac{r_0}{\|\mathbf{y} - \mathbf{x}_0\|_{\Sigma_0}} \in (0, 1),$$

due to  $\mathcal{O}$  being a sphere under a  $\Sigma_0$  weighted norm. The point  $\mathbf{x}^* := \mathbf{x}(r^*)$  lies on the boundary.

$$\Rightarrow \mathbf{y} - \mathbf{x}^* = (1 - r^*) \cdot (\mathbf{y} - \mathbf{x}_0).$$

Since the closest point along an arbitrary ray is never closer than the true minimum, we have

$$\begin{aligned} D_{\Sigma}^{\min}(\mathbf{y}, \mathcal{O}) &= \min_{\mathbf{x} \in \mathcal{O}} \|\mathbf{y} - \mathbf{x}\|_{\Sigma}, \\ &\leq \|\mathbf{y} - \mathbf{x}^*\|_{\Sigma}, \\ &= (1 - r^*) \cdot \|\mathbf{y} - \mathbf{x}_0\|_{\Sigma}, \\ &= \left(1 - \frac{r_0}{\|\mathbf{y} - \mathbf{x}_0\|_{\Sigma_0}}\right) \|\mathbf{y} - \mathbf{x}_0\|_{\Sigma}. \end{aligned}$$

**Condition for equality.** The condition for  $\mathbf{y} \in \mathcal{O}$  case is trivial. For  $\mathbf{y} \notin \mathcal{O}$ , the minimum distance from  $\mathbf{y}$  to the ellipsoid occurs at the boundary intersection of the ray  $\mathbf{x}(r)$  iff the gradient of the *objective*  $\Sigma(\mathbf{y} - \mathbf{x})$  is parallel to the gradient of the *constraint*  $\Sigma_0(\mathbf{x} - \mathbf{x}_0)$  at the intersection point. Since  $(\mathbf{x} - \mathbf{x}_0)$  is proportional to  $(\mathbf{y} - \mathbf{x}_0)$  along the ray, this optimality condition requires:

$$\Sigma(\mathbf{y} - \mathbf{x}_0) \propto \Sigma_0(\mathbf{y} - \mathbf{x}_0) \iff (\mathbf{y} - \mathbf{x}_0) \propto \Sigma^{-1} \Sigma_0(\mathbf{y} - \mathbf{x}_0).$$

Thus, the ray bound is exact if and only if  $(\mathbf{y} - \mathbf{x}_0)$  is an (any) eigenvector of  $\Sigma^{-1} \Sigma_0$ .  $\square$



**Lemma 5** (centroid-radius interaction bound via coupling). For  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma_0)$ , consider two constant vectors  $\mathbf{b}_1, \mathbf{b}_2 \in \mathbb{R}^m$  (representing **centroid displacements**) and two symmetric positive definite matrices  $\Sigma_1, \Sigma_2 \succ \mathbf{0}$  (representing respective **radius metrics**). Define the interaction integral:

$$J(\mathbf{b}, \Sigma) := \mathbb{E}_{\mathbf{x}} \left[ |\mathbf{b}^\top \mathbf{x}| \cdot \sqrt{\mathbf{x}^\top \Sigma \mathbf{x}} \right].$$

If  $(\mathbf{b}_2, \Sigma_2)$  has a strictly shorter “whitened” centroid and a strictly narrower radius (PSD-wise), i.e.,

1. **Centroid Contraction:**  $\|\mathbf{b}_2\|_{\Sigma_0} < \|\mathbf{b}_1\|_{\Sigma_0}$ ,
2. **Radius Contraction:**  $\Sigma_2 \prec \Sigma_1$ ,

then the interaction term strictly decreases:

$$J(\mathbf{b}_2, \Sigma_2) < J(\mathbf{b}_1, \Sigma_1).$$

*Proof.* To evaluate the integral, we transform it into spherically symmetric coordinates (whitening).

**Whitening.** We can express the data vector  $\mathbf{x}$  as a linear transformation of a standard normal vector  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_m)$  such that  $\mathbf{x} = \Sigma_0^{1/2} \mathbf{z}$ . Substituting this into the centroid and radius terms:

$$\text{Centroid: } |\mathbf{b}^\top \mathbf{x}| = |\mathbf{b}^\top \Sigma_0^{1/2} \mathbf{z}| = |(\Sigma_0^{1/2} \mathbf{b})^\top \mathbf{z}| = |\tilde{\mathbf{b}}^\top \mathbf{z}|,$$

$$\text{Radius: } \sqrt{\mathbf{x}^\top \Sigma \mathbf{x}} = \sqrt{\mathbf{z}^\top \Sigma_0^{1/2} \Sigma \Sigma_0^{1/2} \mathbf{z}} = \sqrt{\mathbf{z}^\top \tilde{\Sigma} \mathbf{z}},$$

where  $\tilde{\mathbf{b}} := \Sigma_0^{1/2} \mathbf{b}$  is the whitened centroid, and  $\tilde{\Sigma} := \Sigma_0^{1/2} \Sigma \Sigma_0^{1/2}$  is the whitened radius metric.

**Rotational symmetry (coupling).** The expectation is now over the standard normal variable  $\mathbf{z}$ :

$$J = \mathbb{E}_{\mathbf{z}} \left[ |\tilde{\mathbf{b}}^\top \mathbf{z}| \cdot \sqrt{\mathbf{z}^\top \tilde{\Sigma} \mathbf{z}} \right].$$

Since the distribution of  $\mathbf{z}$  is spherically symmetric (invariant to rotations), the distribution of the dot product  $\tilde{\mathbf{b}}^\top \mathbf{z}$  depends only on the length of  $\tilde{\mathbf{b}}$ . We can conceptually rotate the coordinate system for each scenario such that  $\tilde{\mathbf{b}}$  aligns with the first basis vector  $\mathbf{e}_1$ . In this rotated frame,  $|\tilde{\mathbf{b}}^\top \mathbf{z}| = \|\tilde{\mathbf{b}}\| \cdot |z_1|$ . Crucially, note that  $\|\tilde{\mathbf{b}}\| = \|\Sigma_0^{1/2} \mathbf{b}\|_2 = \sqrt{\mathbf{b}^\top \Sigma_0 \mathbf{b}} = \|\mathbf{b}\|_{\Sigma_0}$ . Thus:

$$J(\mathbf{b}, \Sigma) = \|\mathbf{b}\|_{\Sigma_0} \cdot \mathbb{E}_{\mathbf{z}} \left[ |z_1| \cdot \sqrt{\mathbf{z}^\top \tilde{\Sigma} \mathbf{z}} \right].$$

**Comparison.** We now compare  $J_1 = J(\mathbf{b}_1, \Sigma_1)$  and  $J_2 = J(\mathbf{b}_2, \Sigma_2)$ .

$$\begin{aligned} J_2 &= \|\mathbf{b}_2\|_{\Sigma_0} \cdot \mathbb{E}_{\mathbf{z}} \left[ |z_1| \cdot \sqrt{\mathbf{z}^\top \tilde{\Sigma}_2 \mathbf{z}} \right] \\ &< \|\mathbf{b}_1\|_{\Sigma_0} \cdot \mathbb{E}_{\mathbf{z}} \left[ |z_1| \cdot \sqrt{\mathbf{z}^\top \tilde{\Sigma}_2 \mathbf{z}} \right] && \text{(by centroid contraction)} \\ &< \|\mathbf{b}_1\|_{\Sigma_0} \cdot \mathbb{E}_{\mathbf{z}} \left[ |z_1| \cdot \sqrt{\mathbf{z}^\top \tilde{\Sigma}_1 \mathbf{z}} \right] && \text{(by radius contraction)} \\ &= J_1. \end{aligned}$$

The second inequality holds because  $\Sigma_2 \prec \Sigma_1$  implies  $\tilde{\Sigma}_2 \prec \tilde{\Sigma}_1$ , so  $\mathbf{z}^\top \tilde{\Sigma}_2 \mathbf{z} < \mathbf{z}^\top \tilde{\Sigma}_1 \mathbf{z}$  for all  $\mathbf{z} \neq \mathbf{0}$ . Since  $|z_1|$  is non-negative and not always zero, the expectation strictly decreases.  $\square$

**Lemma 6** (sandwich bounds for SPD-PSD weighted norms). *For  $n \times n$  matrices  $A \succ 0$ ,  $B \succcurlyeq 0$ , denote the pseudo-inverse as  $B^\dagger$ , and  $\Pi_B := B^\dagger B$  projects onto  $\text{range}(B)$ . Then, for any  $x \in \mathbb{R}^n$ ,*

$$\underbrace{\frac{\kappa}{1+\kappa}}_{\text{shrinkage factor} \leq 1} \cdot \|\Pi_B x\|_A^2 \leq x^\top A(A+B)^{-1} B x \leq \|\Pi_B x\|_A^2,$$

*for bounded minimum positive eigenvalue  $\kappa := \lambda_{\min}^+(A^{-1}B) < \infty$ . Equality holds for upper bound iff  $x \perp B$ , and lower bound iff  $\Pi_B x$  is entirely in eigen-space of  $A^{-1}B$  corresponding to  $\kappa$ .*

*Proof.* From Lemma 7, we have  $A = S^\top S$  and  $B = S^\top D S$  for invertible  $S$  and diagonal  $D \succcurlyeq 0$ . Note that  $D$  are eigenvalues of  $A^{-1}B$  by cyclic permutation invariance (i.e.,  $\lambda(AB) = \lambda(BA)$ ).

Define the change of basis  $z := Sx$ . Then  $x = S^{-1}z$ , and

$$\begin{aligned} x^\top A(A+B)^{-1} B x &= x^\top S^\top S(S^\top S + S^\top D S)^{-1} S^\top D S x \\ &= x^\top S^\top S(S^\top (I + D) S)^{-1} S^\top D S x \\ &= x^\top S^\top S S^{-1} (I + D)^{-1} S^{-\top} S^\top D S x \\ &= x^\top S^\top (I + D)^{-1} D S x \\ &= z^\top (I + D)^{-1} D z \\ &= \sum_i \frac{D_{ii}}{1 + D_{ii}} z_i^2. \end{aligned}$$

Similarly, for the projected norm, noting that  $\Pi_B = S^{-1} D^\dagger D S$  and  $\|x\|_A^2 = \|Sx\|_2^2$ :

$$\|\Pi_B x\|_A^2 = \|S(S^{-1} D^\dagger D S)x\|_2^2 = \|D^\dagger D z\|_2^2 = \sum_{i: D_{ii} > 0} z_i^2.$$

**Upper bound.** Since  $\frac{D_{ii}}{1+D_{ii}} < 1$  for all  $D_{ii} > 0$ , the following inequality is strict for any  $z_i \neq 0$ .

$$x^\top A(A+B)^{-1} B x = \sum_i \frac{D_{ii}}{1 + D_{ii}} z_i^2 \leq \sum_i z_i^2 = \|\Pi_B x\|_A^2.$$

And equality holds iff  $z_i = 0$  for all active indices, which implies  $\Pi_B x = 0$  (i.e.,  $x \perp B$ ).

**Lower Bound.** The function  $f(d) = \frac{d}{1+d}$  is monotonically increasing for  $d \geq 0$ . Restricting our attention to the support of the vector (indices where  $D_{ii} > 0$ ), we define  $\kappa = \min\{D_{ii} : D_{ii} > 0\}$ . It follows that for every active index,  $\frac{D_{ii}}{1+D_{ii}} \geq \frac{\kappa}{1+\kappa}$ . Summing over the support:

$$\sum_{i: D_{ii} > 0} \frac{D_{ii}}{1 + D_{ii}} z_i^2 \geq \sum_{i: D_{ii} > 0} \frac{\kappa}{1 + \kappa} z_i^2 = \frac{\kappa}{1 + \kappa} \|\Pi_B x\|_A^2.$$

For the inequality to become an equality, we require  $\frac{D_{ii}}{1+D_{ii}} = \frac{\kappa}{1+\kappa}$  for every index  $i$  where  $z_i \neq 0$ . This implies  $D_{ii} = \kappa$  for all contributing dimensions. Geometrically, this means the vector  $x$  (after projection) must align only with the directions associated with the minimum eigenvalue  $\kappa$ .  $\square$

**Lemma 7** (SPD, PSD joint denationalization via congruence Akbar et al. (2025)). *For any  $n \times n$  matrices  $A \succ 0$ ,  $B \succcurlyeq 0$ , there exists an  $n \times n$  invertible  $S$  and non-negative diagonal  $D$  such that*

$$A = S^\top S, \quad B = S^\top D S.$$

*Proof.* See (Akbar et al., 2025, Lem. 2), cf. (Horn & Johnson, 1985, Thm. 7.6.4, p. 465).  $\square$

## USE OF LARGE LANGUAGE MODELS

A large language model (LLM) was utilized as a writing assistant to help refine the prose, improve clarity, and ensure a consistent narrative tone during the preparation of this manuscript. The human authors directed this process, take full responsibility for the final content, and are solely responsible for all scientific contributions of this work.