Symmetry as Intervention; Causal Estimation with Data Augmentation Anonymous Authors

Abstract

To our knowledge, we provide the first analysis of causal estimation under hidden confounding using only observational (X,Y) data and knowledge of symmetries in data generation via data augmentation (DA) transformations. We show that such DA is equivalent to interventions on the treatment X, mitigating bias from hidden confounding, and that framing DA as a relaxation of instrumental variables (IVs)—sources of X randomization that are conditionally independent of the outcome Y—can further improve causal estimation beyond simple DA. **Keywords:** Causal Inference, Intervention, IV Regression, Invariance, Data Augmentation

1. Preliminaries

For treatment $X \in \mathcal{X} \subseteq \mathbb{R}^m$, outcome $Y \in \mathcal{Y} \subseteq \mathbb{R}^l$ in the structural equation model (SEM) \mathfrak{M}

$$X = \tau(Y, Z, C, N_X), \quad Y = f(X) + \epsilon(C) + N_Y, \quad \text{s.t.} \quad \boxed{\xi := Y - f(X), \quad \mathbb{E}[\xi] = 0,} \quad (1)$$

where Z, C, N_X, N_Y are exogenous, we want to estimate $f \in \mathcal{H} := \{h : \mathcal{X} \to \mathcal{Y}\}$ from $\mathbb{P}^{\mathfrak{M}}_{X,Y}$. When $X \perp \!\!\! \perp \xi$, we estimate f via empirical risk minimization (ERM) given a convex loss ℓ ,

$$R_{\mathrm{ERM}}^{\mathfrak{M}}(h) := \mathbb{E}^{\mathfrak{M}}[\ell(Y, h(X))], \qquad \hat{h}_{\mathrm{ERM}}^{\mathfrak{M}} := \operatorname{argmin}_{h \in \mathcal{H}} R_{\mathrm{ERM}}^{\mathfrak{M}}(h).$$
 (2)

For finite n samples $\mathcal{D} := \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=0}^n$, data augmentation (DA) is used to reduce estimation variance (Lyle et al., 2020) via multiple random augmentations $(G\mathbf{x}_i, \mathbf{y}_i)$ per sample in the risk

$$R_{\mathrm{DA}_G + \mathrm{ERM}}^{\mathfrak{M}}(h) := \mathbb{E}^{\mathfrak{M}}[\ell(Y, h(GX))], \qquad G \sim \mathbb{P}_G.$$
 (3)

However, generally $X \not\perp \!\!\! \downarrow \xi$ due to which the ERM minimizer is biased (Pearl, 2009). This bias is known as the *confounding bias* and X, Y are said to be confounded. Confounding is removed via an *intervention* do(X := X') that sets X to some i.i.d. X', now yielding the *causal risk*

$$R_{\mathrm{CR}}^{\mathfrak{M}}(h) := R_{\mathrm{ERM}}^{\mathfrak{M}; \mathrm{do}(X)}(h) = R_{\mathrm{ERM}}^{\mathfrak{M}; \mathrm{do}(X := X')}(h), \qquad \text{s.t.} \qquad X' \sim \mathbb{P}_X^{\mathfrak{M}}. \tag{4}$$

Where do(X) denotes such interventions. Minimizers of Eq. (4) identify f and are robust predictors to $\mathbb{P}_X^{\mathfrak{M}}$ shifts over $\operatorname{supp}(\mathbb{P}_X^{\mathfrak{M}})$ (Christiansen et al., 2022). Define causal excess risk (CER)

$$CER_{\mathfrak{M}}(h) := R_{CR}^{\mathfrak{M}}(h) - R_{CR}^{\mathfrak{M}}(f),$$

to capture estimation error by removing irreducible noise from Eq. (4), so that $CER_{\mathfrak{M}}(f) = 0$. In practice, interventions are often unavailable. A common workaround is to use auxiliary variables. One approach is that of instrumental variable (IV) regression (Belsley, 1988), where an instrument Z satisfies: (i) **treatment relevance** $Z \not\perp\!\!\!\perp X$, (ii) **exclusion** $Z \perp\!\!\!\perp Y^{\mathfrak{M}; do(X)}$, (iii) **un-confoundedness** $Z \perp\!\!\!\perp \xi$, and (iv) **outcome relevance** $Y \not\perp\!\!\!\perp Z$. Now, Eq. (1) gives

$$\mathbb{E}^{\mathfrak{M}}[Y \mid Z] = \mathbb{E}^{\mathfrak{M}}[f(X) \mid Z]. \tag{5}$$

which admits consistent estimation of f and can be solved by minimizing the following risk

$$R_{\text{IV}}^{\mathfrak{M}}(h) := \mathbb{E}^{\mathfrak{M}} \Big[\ell \Big(Y, \ \mathbb{E}^{\mathfrak{M}}[h(X) | Z] \Big) \Big]. \tag{6}$$

^{1.} Assume all SEMs under discussion entail unique observational distributions $\mathbb{P}_{X,Y}^{\mathfrak{M}}$. Details in Appendix B.

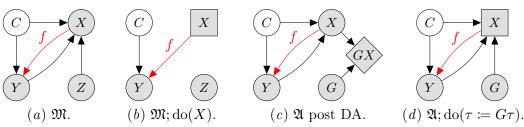


Figure 1: Graphs of respective SEMs; (a) Z is an IV w.r.t. confounded (X,Y). (b) Graph obtained via intervention on X in \mathfrak{M} ; IV regression simulates this intervention with only observational data. (c) Graph for DA. (d) Graph for soft intervention. Observational distributions of (GX,Y,G,C) in (c) and (X,Y,G,C) in (d) are identical.

2. Causal Estimation with Data Augmentation

Problem setup. We discuss the following SEM \mathfrak{A} for exogenous C, N_X, N_Y and $X \not\perp \!\!\! \perp \xi$,

$$X = \tau(Y, C, N_X), \quad Y = f(X) + \epsilon(C) + N_Y, \quad \text{s.t.} \quad \begin{bmatrix} \xi := Y - f(X), & \mathbb{E}[\xi] = 0. \end{bmatrix}$$
 (7)

Consider also a data augmentation with respect to which f is invariant (Chen et al., 2020). The action of a group \mathcal{G} is a mapping $\delta: \mathcal{X} \times \mathcal{G} \to \mathcal{X}$ compatible with the group operation. We write $\mathbf{g}\mathbf{x} := \delta(\mathbf{x}, \mathbf{g})$ as shorthand and say that f is \mathcal{G} -invariant if $f(\mathbf{g}\mathbf{x}) = f(\mathbf{x})$, $\forall (\mathbf{g}, \mathbf{x}) \in \mathcal{G} \times \mathcal{X}$. We refer to such a map $\mathbf{g}\mathbf{x}$, henceforth assumed to be continuous in \mathbf{x} , as a valid outcome-invariant DA transformation parameterized by the vector $\mathbf{g} \in \mathcal{G}$. Let \mathcal{G} have a (unique) normalized Haar measure and $\mathbb{P}_G^{\mathfrak{A}}$ the corresponding distribution defined over it.

The task. Given samples for only $(X,Y) \sim \mathbb{P}_{X,Y}^{\mathfrak{A}}$ and a valid outcome invariant DA parameterized by $G \sim \mathbb{P}_{G}^{\mathfrak{A}}$, we want to improve estimation of f compared to standard ERM.

Now, take a *soft* intervention on $\mathfrak A$ where we replace the mechanism τ of X with $G\tau$. Abusing notation, we represent this SEM by $\mathfrak A$; $\operatorname{do}(\tau \coloneqq G\tau)$, its graph depicted in Fig. 1(d).² Comparing the DA mechanism in $\mathfrak A$ (Fig. 1(c)) and the intervention $\mathfrak A$; $\operatorname{do}(\tau \coloneqq G\tau)$ (Fig. 1(d)):

Observation 1 (soft intervention with DA) $\mathbb{P}^{\mathfrak{A}}_{GX,Y,G,C}$ and $\mathbb{P}^{\mathfrak{A};\operatorname{do}(\tau:=G\tau)}_{X,Y,G,C}$ are identical.

We can hence treat samples generated from $\mathfrak A$ via DA as if they were instead generated from $\mathfrak A$; $\operatorname{do}(\tau \coloneqq G\tau)$ by intervening on X. This allows us to rewrite the risk from Eq. (3) as $R^{\mathfrak A}_{\operatorname{DA}_G+\operatorname{ERM}}(h)=R^{\mathfrak A;\operatorname{do}(\tau \coloneqq G\tau)}_{\operatorname{ERM}}(h)$, to emphasize that DA is equivalent to a (soft) intervention and as such can mitigate confounding bias when estimating f, as shown in the next example.

Example 1 (a linear Gaussian DA example) For $\kappa, \sigma > 0$, non-zero $\Gamma, \mathbf{T} \in \mathbb{R}^{* \times m}$ and $\boldsymbol{\tau}^{\top}, \mathbf{f}, \boldsymbol{\epsilon} \in \mathbb{R}^{m}$ such that $\mathbf{f}^{\top} \boldsymbol{\tau}^{\top} \neq \kappa^{-1}$ so that the following SEM \mathfrak{A} is solvable in $(X, Y)^{3}$

$$X = \kappa \cdot \boldsymbol{\tau}^{\top} Y + \mathbf{T}^{\top} C + \sigma N_X, \qquad Y = \mathbf{f}^{\top} X + \kappa \cdot \boldsymbol{\epsilon}^{\top} C + \sigma N_Y, \qquad GX \coloneqq X + \gamma \cdot \boldsymbol{\Gamma}^{\top} G,$$

where G, C, N_X, N_Y are conformable, centered Gaussian vectors, κ determines how much (X,Y) are confounded and range $(\mathbf{\Gamma}^\top) \subseteq \text{null}(\mathbf{f}^\top)$ to make GX a valid outcome invariant DA.

We evaluate an estimate $\hat{\mathbf{h}}^{\mathcal{D}}$ using CER. For squared loss and covariance $\mathbf{\Sigma}_X^{\mathfrak{A}}$ in Example 1,

$$CER_{\mathfrak{A}}(\hat{\mathbf{h}}^{\mathcal{D}}) = \left\| \hat{\mathbf{h}}^{\mathcal{D}} - \mathbf{f} \right\|_{\Sigma_{X}^{\mathfrak{A}}}^{2}.$$
 (8)

^{2.} For any \mathfrak{A} with unique distribution, \mathfrak{A} ; do($\tau := G\tau$) also has a unique distribution (proof in Appendix H.3).

^{3.} See Appendix B, Lemma 3 for details on solving for and sampling (X,Y) in such linear, simultaneous SEMs.

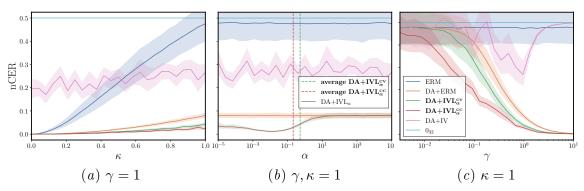


Figure 2: Simulation experiment for the linear Gaussian SEM in Example 1. κ and γ control the amount of confounding and *strength* of DA respectively. α is the IVL regularization parameter. Each data-point averages nCER over 32 trials with a 95% CI.

Theorem 1 (causal estimation with DA+ERM) For SEM 21 in Example 1, we have

$$\operatorname{CER}_{\mathfrak{A}}\left(\hat{\mathbf{h}}_{DA_{G}+ERM}^{\mathfrak{A}}\right) \leq \operatorname{CER}_{\mathfrak{A}}\left(\hat{\mathbf{h}}_{ERM}^{\mathfrak{A}}\right), \quad equality \ iff \quad \mathbb{E}^{\mathfrak{A}}[GX \mid G] \perp \mathbb{E}^{\mathfrak{A}}[X \mid \xi] \quad a.s.$$

Proof See Appendix H.4 for the proof.

That is, DA improves causal estimation iif it targets spurious features of X. Domain knowledge may therefore be needed to design such DA. Still, with outcome invariance, DA is never worse than ERM; allowing regularization at worst, and mitigating confounding bias at best.

We once again point our attention to the graph of \mathfrak{A} ; do($\tau := G\tau$) from Fig. 1(d) to see:

Observation 2 (IV-like DA parameters) In SEM \mathfrak{A} ; do($\tau := G\tau$), the DA parameters G satisfy IV properties (i) through (iii). We refer to such an IV relaxation as IV-like (IVL).

This IV relaxation may render an ill-posed Eq. (5), so we suggest the regularization $R^{\mathfrak{M}}_{\mathrm{IVL}_{\alpha}}(h) \coloneqq R^{\mathfrak{M}}_{\mathrm{IV}}(h) + \alpha R^{\mathfrak{M}}_{\mathrm{ERM}}(h)$ as IVL regression, discussed separately in Appendix E. When composed with DA in \mathfrak{A} now gives $R^{\mathfrak{A}}_{\mathrm{DA}_G+\mathrm{IVL}_{\alpha}}(h) = R^{\mathfrak{A};\mathrm{do}(\tau \coloneqq G\tau)}_{\mathrm{IVL}_{\alpha}}(h)$. The next results follow.

Corollary 1 (worst-case DA with DA+IVL regression) For SEM 31 in Example 1,

$$\hat{\mathbf{h}}_{DA_G+IVL_{\alpha}}^{\mathfrak{A}} \in \operatorname*{argmin}_{\mathbf{h}} \max_{\mathbf{g} \in \mathcal{G}_{\alpha}} R_{DA_{\mathbf{g}}+ERM}^{\mathfrak{A}}(\mathbf{h}), \quad s.t. \quad \mathcal{G}_{\alpha} \coloneqq \bigg\{ \mathbf{g} \, \bigg| \, \boldsymbol{\Gamma}^{\top} \mathbf{g} \mathbf{g}^{\top} \boldsymbol{\Gamma} \preccurlyeq \bigg(\frac{1}{\alpha} + 1 \bigg) \boldsymbol{\Gamma}^{\top} \boldsymbol{\Sigma}_{G}^{\mathfrak{A}} \boldsymbol{\Gamma} \bigg\}.$$

Proof The result follows from Observation 1, Observation 2 and Theorem 2.

Corollary 2 (causal estimation with DA+IVL regression) In Example 1, $\alpha, \gamma < \infty$,

$$\operatorname{CER}_{\mathfrak{A}}\left(\hat{\mathbf{h}}_{DA_{G}+IVL_{\alpha}}^{\mathfrak{A}}\right) \leq \operatorname{CER}_{\mathfrak{A}}\left(\hat{\mathbf{h}}_{DA_{G}+ERM}^{\mathfrak{A}}\right), \quad equality \ iff \quad \mathbb{E}^{\mathfrak{A}}[GX \,|\, G] \perp \mathbb{E}^{\mathfrak{A}}[X \,|\, \xi] \quad a.s.$$

Proof The result follows directly from Theorem 3 and Observation 2.

Using DA parameters as IVL therefore simulates a worst-case, or adversarial application of DA within a set of transforms \mathcal{G}_{α} . Of course Corollary 1 can also be viewed as a predictor that generalizes to treatment interventions encoded by \mathcal{G}_{α} . As is intuitive, such a worst-case intervention improves causal estimation so long as the features of X intervened along include some that are spurious (Corollary 2). DA and IVL regression may therefore be used in composition if the application can benefit from regularization and/ or better prediction generalization across DA-induced interventions, with a "bonus" of lower confounding bias if the DA also augments any spurious features of X. The Appendix covers limitations and related work.

AUTHORS

Extended Abstract Track

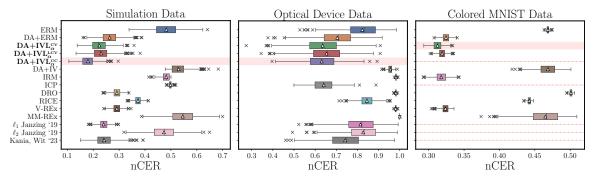


Figure 3: Experiment results; common domain generalisation benchmarks compared against the ERM, DA+ERM and DA+IV baselines, including DA+IVL.

3. Experiments

We empirically evaluate DA's effectiveness in reducing hidden confounding bias in the finite-sample regime. Since our focus is on generalizing across interventions rather than i.i.d. generalization, we fix augmented data size to match the original throughout all experiments.

We compare against standard ERM, DA, IV regression, and re-purposed domain generalization methods including DRO, IRM, ICP, RICE, V-REx, MM-REx, and causal regularization approaches. For methods requiring additional variables, we replace these with DA parameters G (see Appendix G for implementation details and detailed analysis).

For better interpretability of results, we evaluate using normalized CER (nCER): $nCER_{\mathfrak{M}}(h) = \frac{CER_{\mathfrak{M}}(h)}{CER_{\mathfrak{M}}(h) + CER_{\mathfrak{M}}(h_0)} \in [0, 1]$, where h_0 represents null treatment effect. This has the property that nCER = 0 for ground-truth causal solution but 1 under pure confounding.

Simulation experiment. Using the linear SEM from our theory with m = 32, n = 2048 samples across 32 experiments, we find: (1) ERM degrades with increasing confounding κ , (2) DA alone improves performance, (3) DA+IVL achieves best results while DA+IV is unstable. The cross validation approaches of CC, CV and LCV are explained in Appendix E.

Optical device dataset. On 1000 samples across 12 datasets where hidden confounders affect both webcam-captured images and photo-diode readings, DA+ERM improves over ERM, with DA+IVL outperforming other baselines.

Colored MNIST. Where training labels spuriously correlate with color but correlation flips at test, DA via perturbations to hue/brightness helps reduce confounding. DA+ERM provides substantial gains over ERM, with DA+IVL achieving competitive performance

4. Conclusion

We conclude that re-purposing the widely used variance reduction tool of data augmentation (DA) for reducing hidden confounding bias can be effective under outcome invariance. Crucially, it offers a "no-regret" choice for practitioners; improving causal estimation when targeting spurious features, yet performing no worse than the ERM baseline otherwise. Such mitigation of hidden confounding has direct positive implications for the downstream tasks of robust prediction across shifts in $\mathbb{P}_X^{\mathfrak{A}}$ (Reddy et al., 2025), tighter bounds in partial identification (Kilbertus et al., 2020), and more informative sensitivity analyses De Bartolomeis et al. (2024).

CAUSAL ESTIMATION WITH DATA AUGMENTATION

Extended Abstract Track

References

- 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, editors, *Proceedings of the 41st International Conference on Machine Learning*, volume 235, pages 1709–1729. PMLR, 2024.
- David A. Belsley. Two-or three-stage least squares? Computer Science in Economics and Management, 1:21–30, 1988.
- Aharon Ben-Tal, Laurent El Ghaoui, and Arkadi Nemirovski. *Robust Optimization*. Princeton University Press, Princeton, NJ, 2009. doi: 10.1515/9781400831050.
- Andrew Bennett, Nathan Kallus, and Tobias Schnabel. Deep generalized method of moments for instrumental variable analysis. In *Advances in Neural Information Processing Systems* 32. Curran Associates Inc., 2019.
- Dennis S. Bernstein. *Matrix Mathematics: Theory, Facts, and Formulas*. Princeton University Press, second edition, 2009.
- Stephan Bongers, Patrick Forré, Jonas Peters, and Joris M. Mooij. Foundations of structural causal models with cycles and latent variables. *The Annals of Statistics*, 49(5), 2021. doi: 10.1214/21-AOS2064.
- Nate Breznau. Positive returns and equilibrium: Simultaneous feedback between public opinion and social policy. *Policy Studies Journal*, 45, 2016. doi: 10.1111/psj.12171.
- 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.
- Carl F. Christ. The Cowles Commission's contributions to econometrics at Chicago, 1939-1955. Journal of Economic Literature, 32(1):30–59, 1994.
- 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.
- Kevin A. Clarke. The phantom menace: Omitted variable bias in econometric research. Conflict Management and Peace Science, 22(4):341–352, 2005. doi: 10.1080/07388940500339183.
- Bo Dai, Niao He, Yunpeng Pan, Byron Boots, and Le Song. Learning from Conditional Distributions via Dual Embeddings. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54, pages 1458–1467. PMLR, 2017.

Authors

- Hugh Dance and Benjamin Bloem-Reddy. Causal inference with cocycles, 2024. arXiv:2405.13844.
- Piersilvio De Bartolomeis, Javier Abad Martinez, Konstantin Donhauser, and Fanny Yang. Hidden yet quantifiable: A lower bound for confounding strength using randomized trials. In *International Conference on Artificial Intelligence and Statistics*, pages 1045–1053. PMLR, 2024.
- Y. Dubois et al. Lossy compression for lossless prediction. In NeurIPS, 2021.
- Mordecai Ezekiel. The cobweb theorem. The Quarterly Journal of Economics, 52(2):255–280, 1938.
- 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, editors, Advances in Neural Information Processing Systems, volume 36, pages 70638–70653. Curran Associates, Inc., 2023.
- John Fox. Simultaneous equation models and two-stage least squares. *Sociological Methodology*, 10:130–150, 1979.
- William H. Greene. Econometric analysis. Pearson Education India, 2003.
- Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=lQdXeXDoWtI.
- Alastair R. Hall. Generalized method of moments. In *A Companion to Theoretical Econometrics*, pages 230–255. Wiley, 2003.
- Christina Heinze-Deml, Jonas Peters, and Nicolai Meinshausen. Invariant causal prediction for nonlinear models. *Journal of Causal Inference*, 6(2), 2018.
- Tom Heskes. Bias-variance decompositions: the exclusive privilege of bregman divergences, 2025. URL https://arxiv.org/abs/2501.18581.
- Arthur E. Hoerl and Robert W. Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 42(1):80–86, 2000.
- Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, Cambridge, 1985.
- Zhenglin Huang, Xiaoan Bao, Na Zhang, Qingqi Zhang, Xiao Tu, Biao Wu, and Xi Yang. IPMix: Label-preserving data augmentation method for training robust classifiers. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, Advances in Neural Information Processing Systems, volume 36, pages 63660–63673. Curran Associates, Inc., 2023.

Causal Estimation with Data Augmentation

- Antti Hyttinen, Frederick Eberhardt, and Patrik O. Hoyer. Learning linear cyclic causal models with latent variables. *Journal of Machine Learning Research*, 13:3387–3439, 2012.
- Maximilian Ilse, Jakub M. Tomczak, and Patrick Forré. Selecting data augmentation for simulating interventions. In *International Conference on Machine Learning*, pages 4555–4562. PMLR, 2021.
- Dominik Janzing. Causal regularization. In Advances in Neural Information Processing Systems, volume 32, 2019.
- Dominik Janzing and Bernhard Schölkopf. Detecting confounding in multivariate linear models via spectral analysis. *Journal of Causal Inference*, 6(1), 2018.
- John Johnston. Econometric Methods. McGraw-Hill, New York, second edition, 1971.
- 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, pages 20108–20119, 2020.
- Daphne Koller and Nir Friedman. Probabilistic graphical models: principles and techniques. MIT press, 2009.
- 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*, pages 5815–5826. PMLR, 2021.
- Gustavo Lacerda, Peter Spirtes, Joseph Ramsey, and Patrik O. Hoyer. Discovering cyclic causal models by independent components analysis. In *Proceedings of the Twenty-Fourth Conference on Uncertainty in Artificial Intelligence*, pages 366–374. AUAI Press, 2008.
- Steffen L. Lauritzen and Thomas S. Richardson. Chain graph models and their causal interpretations. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 64(3):321–348, 2002. doi: 10.1111/1467-9868.00340.
- Greg Lewis and Vasilis Syrgkanis. Adversarial generalized method of moments, 2018. arXiv:1803.07164.
- Chaochao Lu, Yuhuai Wu, José Miguel Hernández-Lobato, and Bernhard Schölkopf. Invariant causal representation learning for out-of-distribution generalization. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=-e4EXDWXnSn.
- 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.

AUTHORS

- Divyat Mahajan, Shruti Tople, and Amit Sharma. Domain generalization using causal matching. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pages 7313–7324. PMLR, 2021.
- Margaret Mooney Marini. Women's educational attainment and the timing of entry into parenthood. American Sociological Review, 49(4):491–511, 1984.
- Arash Mastouri, Yuhang Zhu, Limor Gultchin, Anna Korba, Ricardo Silva, Matt Kusner, Arthur Gretton, and Krikamol Muandet. Proximal causal learning with kernels: Two-stage estimation and moment restriction. In *Proceedings of 38th International Conference on Machine Learning (ICML)*, volume 139, pages 7512–7523. PMLR, 2021.
- O. Montasser et al. Transformation-invariant learning and theoretical guarantees for ood generalization. In NeurIPS, 2024.
- Joris M. Mooij, Dominik Janzing, Tom Heskes, and Bernhard Schölkopf. On causal discovery with cyclic additive noise models. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24, 2011.
- John F. Muth. Rational expectations and the theory of price movements. *Econometrica*, 29 (3):315–335, 1961.
- Michael Oberst, Nikolaj Thams, Jonas Peters, and David Sontag. Regularizing towards causal invariance: Linear models with proxies. In Marina Meila and Tong Zhang, editors, Proceedings of the 38th International Conference on Machine Learning, volume 139 of Proceedings of Machine Learning Research, pages 8260–8270. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/oberst21a.html.
- Benjamin I. Page and Calvin C. Jones. Reciprocal effects of policy preferences, party loyalties and the vote. *American Political Science Review*, 73(4):1071–1089, 1979. doi: 10.2307/1953990.
- 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, pages 18143–18156, 2022.

Causal Estimation with Data Augmentation

- Abbavaram Gowtham Reddy, Celia Rubio-Madrigal, Rebekka Burkholz, and Krikamol Muandet. When shift happens confounding is to blame, 2025. URL https://arxiv.org/abs/2505.21422.
- Michael R. Roberts and Toni M. Whited. Chapter 7 endogeneity in empirical corporate finance. In George M. Constantinides, Milton Harris, and Rene M. Stulz, editors, *Handbook of the Economics of Finance*, volume 2 of *Handbook of the Economics of Finance*, pages 493–572. Elsevier, 2013. doi: https://doi.org/10.1016/B978-0-44-453594-8.00007-0.
- D. Romero and S. Lohit. Learning partial equivariances from data. In NeurIPS, 2022.
- 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.
- Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=ryxGuJrFvS.
- H. Shao et al. A theory of pac learnability under transformation invariances. In NeurIPS, 2022.
- Robert Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):267–288, 1996.
- 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, editors, *Advances in Neural Information Processing Systems*, volume 35. Curran Associates, Inc., 2022.
- Ruoyu Wang, Mingyang Yi, Zhitang Chen, and Shengyu Zhu. Out-of-distribution generalization with causal invariant transformations. In 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 375–385, 2022. doi: 10.1109/CVPR52688.2022.00047.
- Chi-Sum Wong and Kenneth S. Law. Testing reciprocal relations by nonrecursive structurale-quation models using cross-sectional data. *Organizational Research Methods*, 2(1):69–87, 1999. doi: 10.1177/109442819921005.
- S. Wong et al. Understanding data augmentation for classification: When to warp? In *DICTA*, 2016.
- Jefrey M. Wooldridge. *Econometric Analysis of Cross Section and Panel Data*. The MIT Press, 2010.
- Feng Xie and David Levinson. How streetcars shaped suburbanization: a Granger causality analysis of land use and transit in the Twin Cities. *Journal of Economic Geography*, 10(3): 453–470, 2010.

AUTHORS

- Liyuan Xu and Arthur Gretton. A neural mean embedding approach for back-door and front-door adjustment, 2022. arXiv:2210.06610.
- 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.
- Arnold Zellner and H. Theil. Three-stage least squares: Simultaneous estimation of simultaneous equations. Econometrica, 30(1):54-78, 1962.
- 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.

Causal Estimation with Data Augmentation

Extended Abstract Track

Appendix Contents

Confounding Bias	13
Simultaneity as Cyclic Structures in Equilibrium	15
Related Work	17
IV Regression	2 0
IV-like Regression	22
Limitations	24
Experiment Supplement	25
G.1 Simulation experiment	25
G.3 Colored-MNIST experiment	28
Proofs	31
H.1 Proof of Proposition 1 – DA induced regularization	31
H.2 Proof of Proposition 2 – IVL regression closed form solution in the linear case	32
$\mathrm{H.3}$ Proof of Proposition 3 – Existence of interventional distribution far a DA	33
H.4 Proof of Theorem 1 – Causal estimation with DA+ERM	35
H.5 Proof of Theorem 2 – Robust prediction with IVL regression	36
H.6 Proof of Theorem 3 – Causal estimation with IVL regression	38
H.7 Miscellaneous supporting lemmas	40
	Simultaneity as Cyclic Structures in Equilibrium Related Work IV Regression IV-like Regression Limitations Experiment Supplement G.1 Simulation experiment G.2 Optical device experiment G.3 Colored-MNIST experiment G.4 Colored-MNIST experiment G.5 Colored-MNIST experiment G.6 Proofs H.1 Proof of Proposition 1 – DA induced regularization H.2 Proof of Proposition 2 – IVL regression closed form solution in the linear case H.3 Proof of Proposition 3 – Existence of interventional distribution far a DA H.4 Proof of Theorem 1 – Causal estimation with DA+ERM H.5 Proof of Theorem 2 – Robust prediction with IVL regression H.6 Proof of Theorem 3 – Causal estimation with IVL regression

Authors

Extended Abstract Track

List of Symbols

The notation is largely borrowed from Peters et al. (2017), with overloading where necessary.

 $\mathbb{R}^{n \times *}$ $n \times *$ Euclidean space; dimension * conformal with & inferred from context.

x Scalar.

 \mathbf{x} Vector. When \mathbf{x}^{\top} is described as a vector, it means \mathbf{x} is a flat $1 \times *$ matrix.

X Matrix.

 \mathcal{X} Set.

X Random vector.

m SEM.

 $X^{\mathfrak{M}}$ Random vector X with its SEM \mathfrak{M} specified when unclear from context.

 $\mathbb{P}_{X}^{\mathfrak{M}}$ Distribution of X entailed by \mathfrak{M} ; superscript dropped if clear from context.

 $\mathbb{E}^{\mathfrak{M}}[X]$ Expectation of X under $\mathbb{P}_X^{\mathfrak{M}}$.

 $\mathbb{V}^{\mathfrak{M}}[X]$ Variance of X under $\mathbb{P}_{X}^{\mathfrak{M}}$.

 $\Sigma_X^{\mathfrak{M}}$ Variance—covariance matrix of X under $\mathbb{P}_X^{\mathfrak{M}}$.

 $\Sigma_{X,Y}^{\mathfrak{M}}$ Cross–covariance matrix of X and Y under $\mathbb{P}_{X,Y}^{\mathfrak{M}}$.

 $do(X := \mathbf{x})$ Intervention — X is explicitly set to \mathbf{x} during data generation.

do(X) Shorthand for do(X := X') where $X' \sim \mathbb{P}_X^{\mathfrak{M}}$ is i.i.d. to X.

 \mathfrak{M} ; do($X := \mathbf{x}$) Intervention SEM.

 $\mathfrak{M}_{X=\mathbf{x}}$ SEM with mechanisms of \mathfrak{M} , but exogenous noise distribution $\mathbb{P}^{\mathfrak{M}}_{N|X=\mathbf{x}}$.

 $\mathfrak{M}_{Y=\mathbf{y}}; \operatorname{do}(X \coloneqq \mathbf{x})$ Counterfactual SEM — intervention SEM of $\mathfrak{M}_{Y=\mathbf{y}}$.

 $X \perp \!\!\!\perp Y$ Random vectors X, Y are independent, i.e. $\mathbb{P}_{Y|X}^{\mathfrak{M}} = \mathbb{P}_{Y}^{\mathfrak{M}}$.

 $\mathbf{x} \perp \mathbf{y} \quad \mathbf{x}, \mathbf{y}$ perpendicular, i.e. $\mathbf{x}^{\top} \mathbf{y} = 0$. For random vectors $X^{\top} Y = 0$ a.s.

 $\hat{h}^{\mathfrak{M}}$ Population (infinite-sample) estimate based on $\mathbb{P}^{\mathfrak{M}}$.

 $\hat{h}^{\mathcal{D}}$ Finite-sample estimate based on dataset \mathcal{D} .

Appendix A. Confounding Bias

Statistical vs. causal inference. The target estimand for the statistical risk in Eq. (2) is the Bayes optimal predictor $\mathbb{E}^{\mathfrak{M}}[Y|X=\cdot]$. Whereas the target estimand for the causal risk in Eq. (4) is the average treatment effect (ATE) Xu and Gretton (2022) defined as

$$f_{\text{ATE}}^{\mathfrak{M}}(\mathbf{x}) \coloneqq \mathbb{E}^{\mathfrak{M}; \text{do}(X \coloneqq \mathbf{x})}[Y \,|\, X = \mathbf{x}].$$

The ATE measures the causal influence of X on Y and is equal to $f(\mathbf{x})$ for the SEM \mathfrak{M} in Eq. (1). As such, statistical inference is concerned with predictions of outcome Y, whereas causal inference is concerned with estimating $f_{\text{ATE}}^{\mathfrak{M}} = f$.

Statistical vs. confounding bias. Both types of inference are subject to bias. Statistical bias arises due to miss-specification of the hypothesis class \mathcal{H} , whereas confounding bias arises due to how the data are generated (making $X \not\perp \xi$). The former is therefore a property of the estimator while the later is a property of the data itself. For an estimator $\hat{h}^{\mathcal{D}}$ with the expected value $\bar{h}(\cdot) := \mathbb{E}_{\mathcal{D}}^{\mathfrak{M}} \left[\hat{h}^{\mathcal{D}}(\cdot) \right]$, these biases are defined as

Statistical bias :=
$$\mathbb{E}^{\mathfrak{M}}[Y|X=\cdot] - \bar{h}(\cdot)$$
,
Confounding bias := $f_{\text{ATE}}^{\mathfrak{M}}(\cdot) - \mathbb{E}^{\mathfrak{M}}[Y|X=\cdot]$,
= $f(\cdot) - \mathbb{E}^{\mathfrak{M}}[Y|X=\cdot]$.

For our model in Eq. (1), confounding bias arises due to (i) the exclusion of the (unobserved) common parent C of X and Y, i.e. a confounder, in the ERM objective (hence fittingly called the *omitted-variable bias* Clarke (2005)) and/or (ii) the model is cyclic so that the noise N_Y may itself correlate with X (called *simultaneity bias* Greene (2003); Fox (1979), or reverse causality Pearl (2009) in the degenerate case). For simplicity we shall refer to as the confounding bias Pearl (2009).

Bias-variance decomposition of the causal risk. Because the treatment X and residual ξ are not correlated under \mathfrak{M} ; do(X) in Eq. (1), for any loss function ℓ that admits a 'clean' or 'additive' bias-variance decomposition Heskes (2025), the causal risk also admits a bias-variance decomposition. Using squared loss, for example, we have for some hypothesis $\hat{h}^{\mathcal{D}}$

$$\Rightarrow R_{\operatorname{CR}}^{\mathfrak{M}}(\hat{h}^{\mathcal{D}})$$

$$= \mathbb{E}^{\mathfrak{M};\operatorname{do}(X)} \left[\left\| Y - \hat{h}^{\mathcal{D}}(X) \right\|^{2} \right],$$

$$= \mathbb{E}^{\mathfrak{M};\operatorname{do}(X)} \left[\left\| f(X) + \xi - \hat{h}^{\mathcal{D}}(X) \right\|^{2} \right], \qquad (Structural eq. of Y.)$$

$$= \mathbb{E}^{\mathfrak{M};\operatorname{do}(X)} \left[\left\| \xi \right\|^{2} \right] + \mathbb{E}^{\mathfrak{M};\operatorname{do}(X)} \left[\left\| f(X) - \hat{h}^{\mathcal{D}}(X) \right\|^{2} \right], \quad (Cross term is 0 as $\xi \perp \!\!\! \perp X^{\mathfrak{M};\operatorname{do}(X)}.)$$$

^{4.} Pearl (Pearl, 2009, p.78,184) similarly uses the term for any bias causing observational vs. interventional deviation; this also aligns with econometrics Roberts and Whited (2013); Greene (2003), where both are lumped together as sources of *endogeneity* (i.e., $X \not\perp\!\!\!\perp \xi$).

Authors

Extended Abstract Track

$$=\underbrace{\mathbb{E}^{\mathfrak{M};\operatorname{do}(X)}\Big[\left\|\xi\right\|^2\Big]}_{\text{irreducible noise}} + \underbrace{\mathbb{E}^{\mathfrak{M}}\Big[\left\|f(X) - \hat{h}^{\mathcal{D}}(X)\right\|^2\Big]}_{\text{estimation error, CER}_{\mathfrak{M}}(\hat{h}^{\mathcal{D}}) =}. \ (\mathbb{P}^{\mathfrak{M}}_{X},\,\mathbb{P}^{\mathfrak{M};\operatorname{do}(X)}_{X} \text{ identical by construction.})$$

We can show by following standard procedure that

$$\mathbb{E}_{\mathcal{D}}^{\mathfrak{M}}\Big[\operatorname{CER}_{\mathfrak{M}}\Big(\hat{h}^{\mathcal{D}}\Big)\Big] = \underbrace{\mathbb{E}_{X}^{\mathfrak{M}}\Big[\left\|f(X) - \bar{h}(X)\right\|^{2}\Big]}_{\operatorname{bias}^{2}} + \underbrace{\mathbb{E}_{\mathcal{D}}^{\mathfrak{M}}\Big[\left\|\bar{h}(X) - \hat{h}^{\mathcal{D}}(X)\right\|^{2}\Big]\Big]}_{\operatorname{variance}}.$$

Since for any population estimate $\hat{h}^{\mathfrak{M}}(X) = \bar{h}(X)$, the CER equals the average (squared) bias in estimation

$$CER_{\mathfrak{M}}(\hat{h}^{\mathfrak{M}}) = \mathbb{E}_{X}^{\mathfrak{M}} \left[\left\| f(X) - \hat{h}^{\mathfrak{M}}(X) \right\|^{2} \right] = \mathbb{E}_{X}^{\mathfrak{M}} \left[\left\| f(X) - \bar{h}(X) \right\|^{2} \right].$$

For a rich enough hypothesis class, the ERM estimate coincides with the Bayes optimal predictor $\hat{h}_{\text{ERM}}^{\mathfrak{M}}(\cdot) = \mathbb{E}^{\mathfrak{M}}[Y|X=\cdot]$ and the CER exactly equals the (average squared) confounding bias as we define it above. For a general estimate $\hat{h}^{\mathcal{D}}$, however, the CER also contains statistical bias. Nevertheless, our claims of "better causal estimation via reducing confounding bias" rest on the fact that we are essentially manipulating the data via DA and/or using treatment randomization sources in the form of IVLs. And recall that confounding bias is a property of the data.

Appendix B. Simultaneity as Cyclic Structures in Equilibrium

Since SEM \mathfrak{M} in Eq. (1) is potentially cyclic, a priori it may entail several or no distributions at all. However, we make the assumption that for all $(\mathbf{x}_0, \mathbf{y}_0) \in \mathcal{X} \times \mathcal{Y}$ the unique limits

$$\mathbf{x} = \lim_{t \to \infty} \mathbf{x}_t = \lim_{t \to \infty} \tau(\mathbf{y}_{t-1}, \mathbf{z}, \mathbf{c}, \mathbf{n}_X), \qquad \mathbf{y} = \lim_{t \to \infty} \mathbf{y}_t = \lim_{t \to \infty} f(\mathbf{x}_{t-1}) + \epsilon(\mathbf{c}) + \mathbf{n}_Y$$

exist for any $(\mathbf{z}, \mathbf{c}, \mathbf{n}_X, \mathbf{n}_Y) \sim \mathbb{P}^{\mathfrak{M}}_{Z,C,N_X,N_Y}$, meaning that the unique distribution entailed by \mathfrak{M} is in this equilibrium state. Of course, if \mathfrak{M} is acyclic, these limits always exist.

Linear cyclic assignments

SEMs with cyclic structures have been well studied both in the linear case by Lauritzen and Richardson (2002); Lacerda et al. (2008); Hyttinen et al. (2012), as well as the non-linear case by Mooij et al. (2011); Bongers et al. (2021). Here we briefly provide a causal interpretation to linear simultaneous equations as SEMs with cyclic assignments.

Consider a square matrix $\mathbf{M} \in \mathbb{R}^{d \times d}$ and the SEM

$$W = \mathbf{M}W + N , \qquad (9)$$

where noise N is exogenous and \mathbf{M} allows for a cyclic structure. We enforce $(\mathbf{I}_d - \mathbf{M})$ to be invertible so that the above equation has a unique solution W for any given N. Re-writing the *structural form* in Eq. (9) into a *reduced form*, the distribution of W is defined by

$$W = (\mathbf{I}_d - \mathbf{M})^{-1} N. \tag{10}$$

One way we can present a causal interpretation of the above solution is to view it as a stationary point to the following sequence of random vectors W_t

$$W_t = \mathbf{M}W_{t-1} + N ,$$

which converges if **M** has a spectral norm strictly smaller than one so that $\mathbf{M}^t \to 0$ as $t \to \infty$. The structural form Eq. (9) essentially describes the iterative application of this operation. And in the limit the distribution of $\lim_{t\to\infty} W^t$ will be the same as the reduced form Eq. (10). Although equivalent, reduced form of a cyclic SEM (if one exists) obscures the causal relations in the data generation process.

Furthermore, we restrict our models to not have any "self-cycles" (an edge from a vertex to itself). So, e.g., the matrix **M** in Eq. (9) has all zero diagonal entries. This simplifies our analysis by providing a simple, intuitive interpretation for our definition of DA in Sec. 2, and also ensures that in the non-linear case the SEM entails a unique, well-defined distribution under mild assumptions Bongers et al. (2021); Lacerda et al. (2008).

Similarly we can write the example SEM $\mathfrak M$ from Example 2 in this (block matrix) form as

$$\underbrace{\begin{bmatrix} X \\ Y \end{bmatrix}}_{W} = \underbrace{\begin{bmatrix} \mathbf{0}_{m \times m} & \boldsymbol{\tau}^{\top} \\ \mathbf{f}^{\top} & \mathbf{0}_{1 \times 1} \end{bmatrix}}_{\mathbf{M}} \underbrace{\begin{bmatrix} X \\ Y \end{bmatrix}}_{W} + \underbrace{\begin{bmatrix} \boldsymbol{\Gamma}^{\top} \\ \mathbf{0}_{1 \times k} \end{bmatrix}}_{X} Z + \begin{bmatrix} \mathbf{T}^{\top} \\ \boldsymbol{\epsilon}^{\top} \end{bmatrix} C + \sigma \cdot \begin{bmatrix} N_{X} \\ N_{Y} \end{bmatrix},$$

For this simple case, $(\mathbf{I}_{(m+1)} - \mathbf{M})$ is always invertible so long as $\mathbf{f}^{\top} \boldsymbol{\tau}^{\top} \neq 1$ from Lemma 3. Or we can also restrict $|\mathbf{f}^{\top} \boldsymbol{\tau}^{\top}| < 1$ to ensure that the spectral norm of \mathbf{M} is strictly smaller than 1. We sample from this SEM by first sampling all of the exogenous variables Z, C, N_X, N_Y and then solving the above system for each sample of X, Y via the reduced form in Lemma 3.

A motivating example

Cyclic SEMs were first discussed in the econometrics literature Christ (1994) to model various observational phenomena, and often solved via 2SLS based IV regression Fox (1979) since it is computationally less costly compared to solving the entire system Belsley (1988). A classic example from economics Ezekiel (1938); Muth (1961) is that of a supply and demand model \mathfrak{M} where the relation of price P of a good with quantity Q of demand can be thought of as a cyclic feed-back loop where producers adjust their price in response to demand of the good and consumers change their demand in response to price of a good. In contrast, a change in consumer tastes or preferences would be an exogenous change on the demand curve and can therefore be used as an IV Z.

consumer demand: $Q = \tau \cdot P + \gamma \cdot Z + N_Q$, producer price: $P = f \cdot Q + N_P$.

Where scalars f, τ are such that $|f \cdot \tau| < 1$ so that the system converges to an equilibrium. We say that the measurements made for P and Q are at the equilibrium state of the market⁵ with zero mean measurement noise N_P, N_Q respectively.

Estimation of causal effects – removing simultaneity bias. If we now want to estimate the effect of demand on price f, standard regression will produce a biased estimate $\hat{f}_{\text{ERM}}^{\mathfrak{M}} = f + \frac{\text{Cov}(Q, N_P)}{\text{Var}(Q)}$ because of the simultaneity causing Q and N_P to be correlated (to see this, substitute model of P into the model of Q). We can now use IV regression to get an unbiased estimate of the effect of demand on price in the market as $\hat{f}_{\text{IV}}^{\mathfrak{M}} = f$.

Robust prediction under treatment interventions – avoiding spurious correlations. Similarly, if the producer wants to predict the effect on demand if price is changed (i.e. intervened on), naive ERM will not be a good choice because it will also capture the spurious correlation from $Q \to P$. We therefore use three-stage-least-squares (3SLS) Zellner and Theil (1962); Belsley (1988) (or similar methods) to estimate the ATE $\hat{\tau}_{3SLS}^{\mathfrak{M}} = \mathbb{E}^{\mathfrak{M}; \text{do}(P:=.)}[Q | P = .]$ where we use the first two stages to estimate $\hat{f}_{IV}^{\mathfrak{M}}$, followed by ERM to regress from the residuals $\hat{N}_P := P - \hat{f}_{IV}^{\mathfrak{M}} \cdot Q$ to Q in the third stage.

Other applications. Cyclic SEMs are commonly used in many disciplines to model reciprocally causal phenomena. Application domains include political science Page and Jones (1979); Breznau (2016), sociology Marini (1984), urban planning and design Xie and Levinson (2010), organizational behavior and psychology Wong and Law (1999), etc.

Lastly, to establish clear relevance to the literature of spurious correlations, we present a novel cyclic SEM interpretation of the popular colored-MNIST task in Appendix G.3, which we argue presents a more intuitive perspective of colored-MNIST as a ATE $\mathbb{E}^{\mathfrak{M}; do(X:=.)}[Y|X=.]$ estimatoin task, which is not immediately obvious in the more familiar DAG perspective.

^{5.} In fact, such feed-back models of supply, demand were initially developed to understand irregular fluctuations of prices/quantities that are observed in some markets when not at equilibrium Ezekiel (1938).

Table 1: Bias-variance analysis of canonical regularization methods compared to outcome invariant transformations. We provide a first analysis of confounding bias in ATE estimation for the later. ↓, ↑, — represent a decrease, increase and no-change in the corresponding metric of interest respectively.

	Type of regularization		
	Outcome invariant transform	Canonical $(\ell_1, \ell_2, \text{ or vanilla DA})$	
Statistical variance	↓ Lyle et al. (2020); Chen et al. (2020)	Chen et al. (2020); Hoerl and Kennard (2000); Tibshirani (1996)	
Statistical bias	Lyle et al. (2020); Chen et al. (2020)	↑ Chen et al. (2020); Hoerl and Kennard (2000); Tibshirani (1996)	
Confounding bias	$\begin{matrix} \downarrow \\ (\text{ours}) \end{matrix}$	↑—↓(causal regularization) Janzing (2019); Kania and Wit (2023); Vankadara et al. (2022)	

Appendix C. Related Work

Domain generalization (DG) methods aim to generalize to unseen test domains via robust optimizatoin (RO) Ben-Tal et al. (2009) over a perturbation set \mathcal{P} of possible test domains $\rho \in \mathcal{P}$ as

$$R_{\mathrm{RO}}^{\mathcal{P}}(h) \coloneqq \max_{\rho \in \mathcal{P}} R_{\mathrm{ERM}}^{\rho}(h),$$

Since generalizing to arbitrary test domains is impossible, the choice of perturbation set encodes one's assumptions about which test domains might be encountered. Instead of making such assumptions a priori, it is often assumed to have access to data from multiple training domains which can inform one's choice of perturbation set. This setting is explored in group distributionally robust optimization (DRO) Sagawa et al. (2020). Variations have been used to mitigate confounding bias and subsequently generalize to treatment interventions when used with interventional data Peters et al. (2016); Dance and Bloem-Reddy (2024), confounder information (i.e. entire graph) Krueger et al. (2021); Huang et al. (2023); Lu et al. (2022) or some proxy thereof in the form of environments Arjovsky et al. (2019). We however, do not assume access to any of these and instead synthesize interventions via DA.

Counterfactual DA strategies have been the primary lens for causal analysis of DA 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). These approaches aim for prediction robustness under treatment interventions and often depend on strong assumptions, such as access to the full SEM Yuan et al. (2024); Feder et al. (2023), auxiliary variables Ilse et al. (2021); Feder et al. (2023); Mahajan et al. (2021); Aloui et al. (2023), or causal graphs Pitis et al. (2022); Armengol Urpí et al. (2024). By contrast, we show that outcome-invariance

AUTHORS

Extended Abstract Track

Table 2: Comparison of our proposed 'outcome invariant DA as a (soft) intervention' framework work with prior works on causal analysis of DA. We argue that other frameworks are less general, requiring access to auxiliary variables, the full graph or treatment mechanisms, all of which are often far less accessible than prior knowledge about symmetries of f. Importantly, our analysis is the first to discuss the effects of such DA simulated interventions on treatment effect estimation.

		Type of DA		
		Outcome invariant (ours)	t Counterfactual	
Target SEM	DA simulates	intervention $\mathfrak{A}; do(\tau := G\tau)$	counterfactual $\mathfrak{A}_{Y=\mathbf{y}}; do(\tau := G\tau)$	
Assumed access	Auxiliary data in addition to (X, Y)	X	back-door (i.e., information about ξ) Ilse et al. (2021); Feder et al. (2023); Mahajan et al. (2021); Aloui et al. (2023)	
	Full graph	X	$\checkmark \\ \mbox{Pitis et al. (2022); Armengol Urpí et al. (2024)}$	
	Structural mechanism	v	treatment mechanism τ Yuan et al. (2024); Feder et al. (2023)	
Analysis	robust prediction across $\mathbb{P}_X^{\mathfrak{A}}$ shifts	./	✓	
	causal estimation of treatment effect (i.e., f)	√	X	

of DA suffices for treatment intervention robustness without invoking *counterfactuals*.⁶ Furthermore, prior works have largely ignored causal effect estimation, often assuming reverse-causal settings where the ATE becomes trivial Ilse et al. (2021); Feder et al. (2023); Yuan et al. (2024). To our knowledge, ours is the first framework to study ATE estimation under DA with minimal structural assumptions. See Tab. 2 for a detailed comparison.

Invariant prediction based methods aim to make predictions based on statistical relationships that remain stable across all domains in \mathcal{P} . A common assumption, for instance, is that $\mathbb{P}_{Y|X}$ is invariant across \mathcal{P} , with only the marginal \mathbb{P}_X allowed to vary. Invariance is also closely linked to causal discovery – under the assumption that causal mechanisms remain stable under interventions on inputs Rothenhäusler et al. (2021). This connection has inspired approaches that enforce invariance conditions to uncover causal structures Peters et al. (2016); Heinze-Deml et al. (2018). IV regression can also be viewed as one such method,

^{6.} Representing an SEM with exogenous noise distribution conditioned on some variable $Y = \mathbf{y}$ by $\mathfrak{A}_{Y=\mathbf{y}}$, the counterfactual SEM $\mathfrak{A}_{Y=\mathbf{y}}$; do($X := \mathbf{x}$) is an intervention do($X := \mathbf{x}$) on this new SEM $\mathfrak{A}_{Y=\mathbf{y}}$. The counterfactual distribution then represents questions like 'After observing $Y = \mathbf{y}$, what would have been had $X = \mathbf{x}$ been true.'

Causal Estimation with Data Augmentation

Extended Abstract Track

where the goal is to learn predictors whose residuals are invariant to the instruments Zhang et al. (2023). More broadly, the principle of invariance, whether motivated by causality or not, has proven useful for improving generalization across heterogeneous settings Rothenhäusler et al. (2021); Arjovsky et al. (2019); Dai et al. (2017).

Causal regularization methods are perhaps the best classification for this work. These aim to design regularization strategies Oberst et al. (2021); Kania and Wit (2023) that reduce confounding bias in order to perform well on the down-stream task on prediction robustness across distribution shifts Reddy et al. (2025). Of these, the most comparable to our work are perhaps those that repurpose canonical regularizers like ℓ_1, ℓ_2 for causal estimation Janzing (2019); Kania and Wit (2023); Vankadara et al. (2022). To the best of our knowledge, we are the first to extend this line of study to DA. Table 1 makes a detailed comparison, including with statistical bias-variance analyses.

Outcome invariant DA as causal regularization. Here we show that outcome invariant DA can have a regularizing effect even in the population case as implied by the following result.

Proposition 1 (DA induced regularization) For SEM \mathfrak{A} from Example 1, given decreasing DA strengths $\gamma_1 > \gamma_2 > 0$, we have

$$\left\| \hat{\mathbf{h}}_{DA_G + ERM}^{\mathfrak{A}} \right| \gamma = \gamma_1 \right\|_{\boldsymbol{\Sigma}_X^{\mathfrak{A}}} \leq \left\| \hat{\mathbf{h}}_{DA_G + ERM}^{\mathfrak{A}} \right| \gamma = \gamma_2 \right\|_{\boldsymbol{\Sigma}_X^{\mathfrak{A}}}, \quad equality \ iff \quad \mathbb{E}^{\mathfrak{A}}[GX \mid G] \perp \mathbb{E}^{\mathfrak{A}}[X \mid \xi].$$

Proof See Appendix H.1 for the proof.

Note that this is fundamentally different from the regularization properties of outcome invariant DA in the un-confounded, finite-sample case as described in Lyle et al. (2020); Chen et al. (2020) in the sense that it reduces confounding bias (Theorem 1) by shrinking the coefficients of $\hat{\mathbf{h}}_{\mathrm{DA}_G+\mathrm{ERM}}^{\mathfrak{A}}$ that correspond to confounded features of X which are augmented by the DA.

Appendix D. IV Regression

Two-stage estimators. Minimizing risk of the form Eq. (6) is known as two-stage IV regression. Another approach for two-stage IV regression is to minimize the risk Mastouri et al. (2021); Rothenhäusler et al. (2021)

$$R_{\text{IVLB}}^{\mathfrak{M}}(h) := \mathbb{E}^{\mathfrak{M}} \left[\left\| \mathbb{E}^{\mathfrak{M}}[Y | Z] - \mathbb{E}^{\mathfrak{M}}[h(X) | Z] \right\|^{2} \right], \tag{11}$$

which can be shown to lower-bound (hence the subscript LB) the surrogate risk in Eq. (6) Mastouri et al. (2021) under squared loss.

$$\Rightarrow R_{\text{IV}}^{\mathfrak{M}}(h)$$

$$= \mathbb{E} \Big[\|Y - \mathbb{E}[h(X)|Z]\|^2 \Big],$$

$$= \mathbb{E} \Big[\|(Y - \mathbb{E}[Y|Z]) + (\mathbb{E}[Y|Z] - \mathbb{E}[h(X)|Z])\|^2 \Big], \text{ (Adding and subtracting } \mathbb{E}[Y|Z].)$$

$$= \mathbb{E} \Big[\|Y - \mathbb{E}[Y|Z]\|^2 \Big] + \mathbb{E} \Big[\|\mathbb{E}[Y|Z] - \mathbb{E}[h(X)|Z]\|^2 \Big]$$

$$+ 2\mathbb{E} \Big[(Y - \mathbb{E}[Y|Z])^{\top} (\mathbb{E}[Y|Z] - \mathbb{E}[h(X)|Z]) \Big],$$

$$= \mathbb{E} \Big[\|Y - \mathbb{E}[Y|Z]\|^2 \Big] + \mathbb{E} \Big[\|\mathbb{E}[Y|Z] - \mathbb{E}[h(X)|Z]\|^2 \Big],$$

$$= \mathbb{E} \Big[\|\mathbb{E}[Y|Z] - \mathbb{E}[h(X)|Z]\|^2 \Big] + \mathbb{E} \Big[\mathbb{E} \Big[(Y - \mathbb{E}[Y|Z])^2 \Big| Z \Big] \Big],$$
(Tower rule and scalar Y.)
$$= \mathbb{E} \Big[\|\mathbb{E}[Y|Z] - \mathbb{E}[h(X)|Z]\|^2 \Big] + \mathbb{E}[\mathbb{V}[Y|Z]] = R_{\text{IV}_{\text{LB}}}^{\mathfrak{M}}(h) + \mathbb{E}[\mathbb{V}[Y|Z]],$$
(13)

where Eq. (13) follows from the definition of conditional variance and we get Eq. (12) by setting the cross term to zero since

$$\Rightarrow \mathbb{E}\left[\left(Y - \mathbb{E}[Y|Z]\right)^{\top} (\mathbb{E}[Y|Z] - \mathbb{E}[h(X)|Z])\right]$$

$$= \mathbb{E}\left[\mathbb{E}\left[\left(Y - \mathbb{E}[Y|Z]\right)^{\top} (\mathbb{E}[Y|Z] - \mathbb{E}[h(X)|Z]) \middle| Z\right]\right], \qquad \text{(Tower rule.)}$$

$$= \mathbb{E}\left[\mathbb{E}\left[\left(Y - \mathbb{E}[Y|Z]\right)^{\top} \middle| Z\right] (\mathbb{E}[Y|Z] - \mathbb{E}[h(X)|Z])\right], \qquad (14)$$

$$= \mathbb{E}\left[\left(\mathbb{E}[Y|Z] - \mathbb{E}[Y|Z]\right)^{\top} (\mathbb{E}[Y|Z] - \mathbb{E}[h(X)|Z])\right], \qquad (14)$$

$$= \mathbb{E}\left[\mathbf{0}^{\top} (\mathbb{E}[Y|Z] - \mathbb{E}[h(X)|Z])\right] = 0,$$

where Eq. (14) follows from the "taking out what is known" rule, i.e.,

$$\mathbb{E}[g(B)A|B] = g(B)\mathbb{E}[A|B]. \tag{15}$$

Generalized method of moments. Another popular approach to solve the IV estimation problem is the *generalized methods of moments (GMM)* Hall (2003); Bennett et al. (2019); Lewis and Syrgkanis (2018) or equivalently the *conditional moment restriction*

(CMR) Mastouri et al. (2021) framework which tries to directly solve for the fact that in Eq. (1) with scalar Y

$$\mathbb{E}^{\mathfrak{M}}[\xi \,|\, Z] = \mathbb{E}^{\mathfrak{M}}[Y - f(X) \,|\, Z] = 0,\tag{16}$$

which holds as a direct consequence of the unconfoundedness property of IV Z, however it is a much weaker assumption on it's own⁷. Equation Eq. (16) implies that for any $q: \mathcal{Z} \to \mathbb{R}$, it holds that

$$\mathbb{E}^{\mathfrak{M}}[(Y - f(X)) \cdot q(Z)] = 0.$$

The GMM-IV estimate of f therefore tries to enforce this condition Hall (2003); Bennett et al. (2019); Lewis and Syrgkanis (2018) by minimizing the risk

$$R_{\text{IV}_{\text{GMM}}}^{\mathfrak{M}}(h) := \sum_{i=1}^{\mu} \mathbb{E}^{\mathfrak{M}} \left[\left(Y - h(X) \right) \cdot q_i(Z) \right]^2 = \left\| \mathbb{E}^{\mathfrak{M}} \left[\left(Y - h(X) \right) \cdot \mathbf{q}(Z) \right] \right\|^2,$$

where $\mathbf{q}(\cdot) \in \mathbb{R}^{\mu}$ represents a vector form of the set of μ arbitrary real-valued functions q_i . A more general form of the above GMM based IV risk is to weight the norm by some SPD **W** Johnston (1971); Hall (2003); Bennett et al. (2019)

$$R_{\text{IV}_{\text{GMM-}\mathbf{W}}}^{\mathfrak{M}}(h) := \left\| \mathbb{E}^{\mathfrak{M}}[(Y - h(X)) \cdot \mathbf{q}(Z)] \right\|_{\mathbf{W}}^{2},$$

which gives the most statistically efficient estimator, minimizing the asymptotic variance, for $\mathbf{W} = \mathbf{\Sigma}_Z^{\mathfrak{M}-1}$ Johnston (1971); Hall (2003); Bennett et al. (2019). We use the same for our non-linear experiments, together with the identity function $\mathbf{q}(Z) = Z$. This gives us the final loss of the form

$$R^{\mathfrak{M}}_{\mathrm{IV}_{\mathrm{GMM}-\boldsymbol{\Sigma}_{Z}^{-1}}}(h) = \left\| \mathbb{E}^{\mathfrak{M}}[Z \cdot (Y - h(X))] \right\|_{\boldsymbol{\Sigma}_{Z}^{-1}}^{2}.$$

And the empirical version of which can be written as follows

$$R_{\text{IV}_{\text{GMM}-\mathbf{\Sigma}_{Z}^{-1}}}^{\mathcal{D}}(h) := \left(\hat{\mathbf{y}} - \mathbf{h}(\hat{\mathbf{X}})\right)^{\top} \hat{\mathbf{Z}} \hat{\mathbf{Z}}^{\dagger} \left(\hat{\mathbf{y}} - \mathbf{h}(\hat{\mathbf{X}})\right), \tag{17}$$

where for dataset samples $(\mathbf{x}_i, y_i, \mathbf{z}_i) \in \mathcal{D}$, we construct the vector $\hat{\mathbf{y}} := [y_0, \dots, y_n]^\top$, matrices $\hat{\mathbf{X}} := [\mathbf{x}_0^\top, \dots, \mathbf{x}_n^\top]^\top$, $\hat{\mathbf{Z}} := [\mathbf{z}_0 \dots \mathbf{z}_n]^\top$ with pseudo-inverse $\hat{\mathbf{Z}}^\dagger$ and define $\mathbf{h}(\hat{\mathbf{X}}) := [h(\mathbf{x}_0), \dots, h(\mathbf{x}_n)]^\top$.

^{7.} Therefore an invalid instrument that does not satisfy the unconfoundedness property, but still satisfies Eq. (16) can also be used here.

Appendix E. IV-like Regression

Faithfulness and outcome-relevance in IVs

Consider the SEM \mathfrak{M} from Sec. 1. The distribution $\mathbb{P}^{\mathfrak{M}}_{X,Y,Z,C}$ is said to be *faithful* to the graph of \mathfrak{M} if it only exhibits independences implied by the graph Peters et al. (2017); Koller and Friedman (2009).⁸ This standard assumption in IV settings renders outcome-relevance implicit and therefore rarely mentioned. In this section we discuss the case where only the first three IV properties are satisfied, i.e. outcome-relevance may not hold. Since such a Z may not be a valid IV, therefore identifiability of f is not possible in general as the problem in Eq. (5) can now be misspecified, having multiple, potentially infinitely many solutions when $Y \perp \!\!\! \perp Z$. Nevertheless, we shall refer to such a Z as IV-like (IVL) to emphasize that while Z may not be an IV, it may still be "instrumental" for reducing confounding bias when estimating f compared to the standard ERM baseline.

ERM regularized IV regression. Despite problem miss-specification for a IVL Z, the target function f remains a minimizer for the IV risk in Eq. (6). Albeit, potentially not unique – for example, a linear h with squared loss leads to an under-determined problem in Eq. (6). We therefore propose a regularized version of the IV risk for such an IVL setting,

$$R_{\text{IVL}_{\alpha}}^{\mathfrak{M}}(h) := R_{\text{IV}}^{\mathfrak{M}}(h) + \alpha R_{\text{ERM}}^{\mathfrak{M}}(h),$$
 (18)

where $\alpha > 0$ is the regularization parameter. The ERM risk as a penalty allows our estimations to have good predictive performance while the IV risk encourages solution search within a subspace where we know f to be present. We refer to minimising Eq. (18) as IVL regression.

Note that the motivation behind IVL regression is not the identifiability of f, but rather potentially better estimation of f by reducing confounding bias. We provide an example.

Example 2 (a linear Gaussian IVL example) For $\sigma > 0$, non-zero $\Gamma, \mathbf{T} \in \mathbb{R}^{* \times m}$ and $\boldsymbol{\tau}^{\top}, \mathbf{f}, \boldsymbol{\epsilon} \in \mathbb{R}^{m}$ such that $\mathbf{f}^{\top} \boldsymbol{\tau}^{\top} \neq 1$ so that the following SEM \mathfrak{M} is solvable in (X, Y)

$$X = \boldsymbol{\tau}^{\top} Y + \boldsymbol{\Gamma}^{\top} Z + \boldsymbol{\mathrm{T}}^{\top} C + \sigma N_X, \qquad Y = \boldsymbol{\mathrm{f}}^{\top} X + \boldsymbol{\epsilon}^{\top} C + \sigma N_Y,$$

where Z, C, N_X, N_Y are conformable, centered Gaussian random vectors and Z is IVL.

Theorem 2 (robust prediction with IVL regression) For SEM M in Example 2,

$$\hat{\mathbf{h}}_{IVL_{\alpha}}^{\mathfrak{M}} \in \underset{\boldsymbol{\zeta} \in \mathcal{P}_{\alpha}}{\operatorname{argmin}} \max_{\boldsymbol{\zeta} \in \mathcal{P}_{\alpha}} R_{ERM}^{\mathfrak{M}; \operatorname{do}\left(\boldsymbol{\Gamma}^{\top}(\cdot) := \boldsymbol{\zeta}\right)}(\mathbf{h}), \quad s.t. \quad \mathcal{P}_{\alpha} := \left\{ \boldsymbol{\zeta} \,\middle|\, \boldsymbol{\zeta} \boldsymbol{\zeta}^{\top} \preccurlyeq \left(\frac{1}{\alpha} + 1\right) \boldsymbol{\Gamma}^{\top} \boldsymbol{\Sigma}_{Z}^{\mathfrak{M}} \boldsymbol{\Gamma} \right\}.$$

Proof See Appendix H.5 for the proof.

Theorem 3 (causal estimation with IVL regression) In Example 2, for $\alpha < \infty$,

$$\operatorname{CER}_{\mathfrak{M}}\left(\hat{\mathbf{h}}_{IVL_{\alpha}}^{\mathfrak{M}}\right) \leq \operatorname{CER}_{\mathfrak{M}}\left(\hat{\mathbf{h}}_{ERM}^{\mathfrak{M}}\right), \quad equality \ iff \quad \mathbb{E}^{\mathfrak{M}}[X|Z] \perp \mathbb{E}^{\mathfrak{M}}[X|\xi] \quad a.s$$

Proof See Appendix H.6 for the proof.

^{8.} Also known as *stability* in some texts (Pearl, 2009, p. 48).

Theorem 2 shows that IVL regression achieves optimal predictive performance across treatment interventions within the perturbation set \mathcal{P}_{α} defined by α . Theorem 3 further states that this strictly reduces confounding bias in estimation of f iff the perturbations align with spurious features of X, as indicated by the equality condition (also necessary for identifiability in linear IV settings Wooldridge (2010); Christiansen et al. (2022)).

Closed form solution in the linear case

The next result gives lets us compute a closed-form solution to the IVL_{α} regression problem in the linear Gaussian case. An empirical version of this is used for our linear experiments.

Proposition 2 (IVL $_{\alpha}$ closed form solution) For SEM \mathfrak{M} in Example 2, $\hat{\mathbf{h}}_{IVL_{\alpha}}^{\mathfrak{M}}$ is the closed form linear OLS solution between

$$X' := aX + b\mathbb{E}[X \mid Z],$$
 $Y' := aY + b\mathbb{E}[Y \mid Z],$

where

$$a := \sqrt{\alpha},$$
 $b := \sqrt{1+\alpha} - \sqrt{\alpha}.$

Proof See Appendix H.2 for the proof.

For the empirical version of Proposition 2 we fit a closed-form OLS regressor between

$$X' := \sqrt{\alpha}X + (\sqrt{1+\alpha} - \sqrt{\alpha})\hat{\mathbf{Z}}\hat{\mathbf{Z}}^{\dagger}X, \qquad Y' := \sqrt{\alpha}Y + (\sqrt{1+\alpha} - \sqrt{\alpha})\hat{\mathbf{Z}}\hat{\mathbf{Z}}^{\dagger}Y.$$

Choice of regularization parameter

Selecting the IVL regularization parameter α in the finite sample setting is not very straightforward. We explore a the approaches that are described below which seem to work well in practice, however some of these may not seem as well motivated since the task at hand is OOD generalization and α is being set via cross-validation with-in the same distribution.

Cross validation (CV), or any variation thereof. We specifically use the following two in our experiments; (i) vanilla CV with 20% samples held-out for validation (ii) level cross validation (LCV) for when Z is discrete, where hold-out data corresponding to 20% of the levels of Z for validation.

Confounder correction (CC), where in a linear setting we follow an approach similar to Janzing (2019) by estimating the length of the true solution f from the observational data \mathcal{D} . We then chose α such that the length of $\hat{h}_{\mathrm{DA+IVL}_{\alpha}}^{\mathcal{D}}$ is closest to the estimated length of the ground truth solution.

AUTHORS

Extended Abstract Track

Appendix F. Limitations

Necessity and practicality of prior knowledge. As discussed in Sec. 2, lower confounding bias is not a 'free lunch' with outcome invariant DA and practitioners may need domain knowledge to construct DA that targets spurious features. Nevertheless, under outcome invariance our methods should at least not perform worse than ERM.

Furthermore, causal estimation from only observational data (X,Y) is generally not possible without untestable assumptions such as the ones me make above. For example, the IV (or IVL) properties of un-confoundedness and exclusion are untestable and have to be justified by domain knowledge. When contextualized in the framework of IVLs (Observation 2), we argue that our assumptions on DA are actually quite practical given that a symmetry-based DA model has precedence in the DA and invariance literature (Chen et al., 2020; Lyle et al., 2020; Shao et al., 2022; Fawzi and Frossard, 2015; Dubois et al., 2021; Petrache and Trivedi, 2023; Montasser et al., 2024; Romero and Lohit, 2022; Zhu et al., 2021; Wong et al., 2016).

$$\underbrace{\underbrace{\text{outcome-invariance} + \text{ spurious targets}}_{\text{popular model for DA}} + \underbrace{\text{spurious targets}}_{\text{benign failure if violated}} \iff \underbrace{\text{un-confoundedness} + \text{exclusion}}_{\text{un-confoundedness}}$$

Therefore, the assumptions required by our methods can be quite practical in many settings where valid IVs (or other auxiliary variables) are scarce, but plausible outcome-invariances (i.e., data augmentations) are abundant.

Choice of α in the finite-sample case. Selecting the IVL regularization parameter in the finite sample setting is not very straightforward. We have outlined several approaches in Appendix E that seem to work well in practice, however some of these may not seem as well motivated since the task at hand is OOD generalization and α is being set via cross-validation with-in the same distribution. Nevertheless, this limitations is not unique to our IVL method and is a general problem is most domain generalization methods Gulrajani and Lopez-Paz (2021).

Appendix G. Experiment Supplement

We began by presenting results in the infinite-sample setting to emphasize that mitigating confounding bias is fundamentally not a sample size issue, i.e., not solvable through traditional regularization alone. In this section, we turn to the finite-sample regime and empirically evaluate the effectiveness of DA in reducing hidden confounding bias. Importantly, we do not use DA for its conventional purpose of augmenting data to improve i.i.d. generalization. Since our focus is on generalizing across interventions, we fix the number of samples in the augmented dataset to match that of the original dataset throughout all experiments.

Finding baselines for evaluating our results is however a challenge – reducing the bias due to hidden confounding in regression estimates having only access to the treatment X and outcome Y is a non-trivial problem. Nevertheless, for the sake of completeness we make an effort to re-purpose existing methods from domain-generalization, invariance learning and causal inference literature to be used as baselines. These methods often require access to additional variables (e.g. IVs, confounders, domains/environments, etc.), and to maintain fairness we will replace these with DA parameters G. Such a comparison is conceptually valid since by virtue of being DG methods, they are essentially solving a robust loss of a similar form as in Corollary 1, giving us meaningful baselines for DA+IVL.

In addition to standard ERM, DA and IV regression, our baselines include DRO Sagawa et al. (2020), invariant risk minimization (IRM) Arjovsky et al. (2019), invariant causal prediction (ICP) Peters et al. (2016), regularization with invariance on causal essential set (RICE) Wang et al. (2022), variance risk extrapolation (V-REx) and minimax risk extrapolation (MM-REx) Krueger et al. (2021). We also compare against causal regularization methods, including Kania and Wit Kania and Wit (2023) and the ℓ_1, ℓ_2 approaches by Janzing Janzing (2019). We discretise G if the method accepts only discrete variables. For IVL regression, we select the regularization parameter α in a variety of ways, including vanilla cross validation (CV), level-based cross validation (LCV) and confounder correction (CC) as described in Appendix E. Other implementation details are provided in Appendix G.

For the methods that use stochastic gradient descent (SGD), we use a learning rate of 0.01, batch size of 256 for 16 epochs. For baselines that require a discrete domains/environments, we uniformly discretise each dimension of G into 2 bins. Higher discretisation bins renders most baselines ineffective since each domain/environment rarely has more than 1 sample. To keep the comparison fair, however, we also discretize G for IVL $_{\alpha}$ regression when using LCV. For the colored MNIST experiment, all CV implementations including baselines use 5-folds for a random search over an exponentially distributed regularization parameter with rate parameter of 1. Same is the case for simulation and optical device experiments, except that DA+IVL methods use a log-uniform distributed regularization parameter over [10⁻⁴, 1]. Since RICE Wang et al. (2022) grows the dataset size by augmenting each sample T times, we provide it a 1/T sub-sample of the original data for fair comparison.

G.1. Simulation experiment

For the finite sample results of the linear SEM \mathfrak{A} from Example 1, by taking m=32, k=31 (dimension of G), $\sigma=0.1$ and fixing $\boldsymbol{\tau}^{\top}=\mathbf{0}$, we sample a new $\mathbf{f}, \boldsymbol{\epsilon}$ and $\mathbf{T} \in \mathbb{R}^{m \times m}$ from a standard normal distribution for each of the 32 experiments for every combination of κ and γ . Each time we construct a $\mathbf{\Gamma} := \mathbf{V}_0$ with k rows as orthonormal basis of null(\mathbf{f}), such that

the SVD of \mathbf{f} is

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

Although this construction of Γ relies on direct knowledge of \mathbf{f} (which is unavailable in practice), we include it here purely for illustrative purposes. We treat access to Γ as our prior structural knowledge about the invariance properties of \mathbf{f} , noting that this information alone is insufficient to recover \mathbf{f} .

We then generate n=2048 samples of (X,Y) for each experiment. For ERM we use a closed form linear OLS solution, for DA+IV, we make use of linear 2SLS. Finally, DA+IVL $_{\alpha}$ was implemented using a closed form linear OLS solution between empirical versions (see Proposition 2) of

$$X' \coloneqq \sqrt{\alpha}X + \left(\sqrt{1+\alpha} - \sqrt{\alpha}\right)\mathbb{E}[X\,|\,Z], \qquad Y' \coloneqq \sqrt{\alpha}Y + \left(\sqrt{1+\alpha} - \sqrt{\alpha}\right)\mathbb{E}[Y\,|\,Z].$$

Our first experimental result in Fig. 2(a) compares the different estimation methods across varying levels of confounding $\kappa \in [0,1]$. As expected, ERM performance degrades with increasing confounding. Applying DA alone already brings us closer to the causal solution, while DA+IVL achieves even better performance. DA+IV regression is unstable and generally performs poorly as it is under-determined.

In the second experiment (Fig. 2(b)), we fix the confounding and DA strengths at $\kappa = \gamma = 1$, and sweep over the regularization parameter $\alpha \in [10^{-5}, 10^5]$ for DA+IVL $_{\alpha}$. The results show that optimal performance is achieved for intermediate values of α , confirming that arbitrarily small values of α , while beneficial in the population setting (as suggested by Theorem 3), are suboptimal in finite samples. We also find that both CV and CC strategies effectively select reasonable values of α .

Finally, we examine sensitivity to the DA strength $\gamma \in [10^{-2.5}, 10]$, fixing $\kappa = 1$. As expected, stronger DA results in stronger interventions on X, which improves causal effect estimation. However, we also observe diminishing returns; when the variation induced by DA is either too small or too large, $\mathrm{DA}+\mathrm{IVL}_{\alpha}$ does not yield significant improvements over the DA+ERM baseline.

For completeness, we also benchmark our approach against other baseline methods on 16 distinct simulation SEMs with 2048 samples each. Aggregated results are presented in Fig. 3 (left most).

For the parameter sweep experiments of Fig. 2, we generate a treatment of dimension m=32, but for the OOD baseline comparison experiment in Fig. 3 we use m=16. Furthermore, for the OOD baseline comparison experiment in Fig. 3, we randomly pick each basis of null(\mathbf{f}) with a probability 2/3 to construct Γ (i.e., we know only some, but not all symmetries of \mathbf{f}).

G.2. Optical device experiment

The dataset from Janzing and Schölkopf (2018) consists of 3×3 pixel images X displayed on a laptop screen that cause voltage readings Y across a photo-diode. A hidden confounder

^{9.} We conjecture that this is due to outcome invariance not holding exactly in practice.

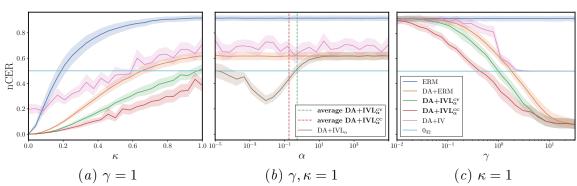


Figure 4: Simulation of the linear Gaussian SEM of Example 1 with the same setting as Fig. 2, but τ^{\top} , **f** sampled uniformly over a unit sphere, representing a cyclic structure.

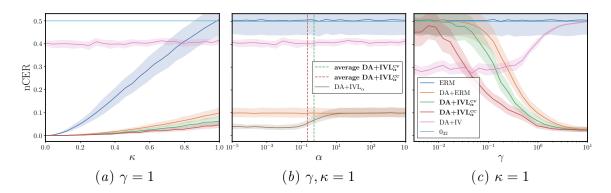


Figure 5: Same experiment as Fig. 2, but with Γ constructed by randomly selecting each basis of null(\mathbf{f}^{\top}) with a probability of 2/3, so that we can simulate the effect of knowing only some symmetries of \mathbf{f} .

C controls two LEDs; one affects the webcam capturing X, the other affects the photo-diode measuring Y. The ground-truth predictor \mathbf{f} is computed by first regressing Y on $(\phi(X),C)$, where $\phi(X)$ are polynomial features of X with degree $d \in \{1,\cdots,5\}$ that best explains the data. The component corresponding to C is then removed to recover \mathbf{f} . We add Gaussian noise $G \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}_X/10)$ for DA and evaluate methods from $\mathbf{??}$ on n=1000 samples across 12 datasets. Figure 3 (middle) shows that DA+ERM improves over ERM, and DA+IVL performs even better, outperforming other baselines.

In this experiment we fit a linear function $h(.) := \mathbf{h} \in \mathbb{R}^m$ for a squared loss in all of our risk metrics. For IVL_{α} regression, we use the closed-form OLS solution from Appendix E. We also use a closed-form solution for ERM, DA+ERM and DA+IV (2SLS) baselines. The rest of the baselines (other than ICP) use SGD.

Most of the datasets in the optical device dataset were best explained by polynomial features of degree 2. We use the same ground-truth degree to fit each of the methods listed in Fig. 3. This is important so as to avoid statistical bias from model miss-specification as our analysis squarely focuses on confounding bias.

G.3. Colored-MNIST experiment

We evaluate on Colored MNIST Arjovsky et al. (2019), where labels are spuriously correlated with image color during training, but this correlation is flipped at test time. We use the same neural architecture and parameters as Arjovsky et al. (2019) across all baselines, training with the IV-based objective described in the Appendix D. DA is implemented via small perturbations to hue, brightness, contrast, saturation, and translation, each parameterized by $G \sim \beta(2,2)$. Although these do not directly manipulate color, the actual spurious feature, they still help reduce confounding. Results in Fig. 3 (rightmost) show that ERM underperforms, DA+ERM provides substantial gains, and DA+IVL $_{\alpha}$ performs competitively with the best DG baselines, with DA+IVL $_{\alpha}^{\rm CV}$ achieving the best overall performance.

We use the same 3-layer neural network (NN) architecture for h across all methods comprising of a fully-connected input layer of input dimension m, hidden layer of input/output dimension 256 and output classification layer with a Sigmoid function. Each layer is separated by an intermediary rectified linear unit activation function. For the IV risk, we use the empirical version of the GMM based risk from Eq. (17).

COLORED-MNIST AS A CYCLIC SEM – FROM INVARIANT PREDICTION TO ESTIMATING CAUSAL EFFECTS

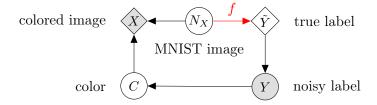


Figure 6: The data generation DAG for colored-MNIST as discussed by the original authors Arjovsky et al. (2019). They aim to learn a predictor $h: \mathcal{X} \to \mathcal{Y}$ such that it is invariant to changes in $\mathbb{P}_{X|Y}$. We argue that this DAG view of colored-MNIST does not make it obvious how the true labeling function $f(\mathbf{x})$ is related to the ATE $\mathbb{E}^{\mathfrak{M};do(X:=\mathbf{x})}[Y|X=\mathbf{x}]$, which we believe is because it is virtually equivalent to the reduced form of our structural form presented in Fig. 7.

In this section we give a cyclic SEM perspective of the colored-MNIST experiment from Arjovsky et al. (2019). The task is binary classification of colored images X from the MNIST dataset into low digits (y = 0 for digits from 0 to 4) and high digits (y = 1 for digits from 5 to 9). The difficulty of the task arises from there being a higher spurious correlation between the color C of the images (c = 0 for blue and c = 1 for green) and (noisy) labels Y as compared to the correlation between the digits in the image and the label.

Consider the following cyclic SEM in Fig. 7.

 $\mathbf{n}_{X} \sim \mathbb{P}_{N_{X}}, n_{Y} \sim \mathbb{B}(0.25), n_{c} \sim \mathbb{B}(e)$ sample all exogenous variables $X = \operatorname{colour}(C, \mathbf{n}_{X})$ apply color C to the image $\tilde{Y} = f(X)$ generate ground-truth label with true labeling function $Y = \operatorname{xor}(\tilde{Y}, n_{Y})$ flip the label with probability 0.25 $C = \operatorname{xor}(Y, n_{C})$ generate color by flipping Y with probability e,

where we first randomly sample an un-colored MNIST image \mathbf{n}_X , and some Bernoulli distributed label noise $n_Y \sim \mathbb{B}(0.25)$ and color noise $n_C \sim \mathbb{B}(e)$ which is different for each environment $e \in \{0.1, 0.2\}$. Then for some initial arbitrary values \mathbf{x}_0 , \tilde{y}_0 , y_0 and c_0 respectively for the observed colored image X, the ground-truth label \tilde{Y} , the observed noisy label Y and the image color C, we iteratively apply the following assignments from the SEM

 $\mathbf{x}_t = \operatorname{colour}(c_{t-1}, \mathbf{n}_X)$ apply color C to the image $\tilde{y}_t = f(\mathbf{x}_{t-1})$ generate ground-truth label with true labeling function $y_t = \operatorname{xor}(\tilde{y}_{t-1}, n_Y)$ flip the label with probability 0.25 $c_t = \operatorname{xor}(y_{t-1}, n_C)$ generate color by flipping Y with probability e,

until they converge while keeping all sampled exogenous variables \mathbf{n}_X, n_Y, n_C fixed. It is straightforward to show that this SEM will converge after a maximum of t = 5 iterations¹⁰ due to the invariance of f to the color of the image C. Furthermore, this stationary-point will be uniquely determined by our exogenous samples \mathbf{n}_X, n_Y, n_C . And this is how we generate one sample (\mathbf{x}, y) for our colored-MNIST experiment. We repeat this process to generate a sample (\mathbf{x}, y) for each of n samples \mathbf{n}_X, n_Y, n_C .

Note that the ground-truth labeling function f can only correctly predict the labels 75% of the time. At test time we flip the correlation between the label Y and the image color C by setting e = 0.9. Also, the above cyclic SEM for colored-MNIST produces the same distribution for (X,Y) as Arjovsky et al. (2019).

The above cyclic SEM perspective of colored-MNIST is interesting because it makes it clear that colored-MNIST is essentially a causal effect estimation task. Specifically, we can estimate the true labeling function f by estimating the ATE $\mathbb{E}^{\mathfrak{M}; do(X:=\mathbf{x})}[Y|X=\mathbf{x}]$ since

$$\mathbb{E}^{\mathfrak{M};\operatorname{do}(X:=\mathbf{x})}[Y \mid X = \mathbf{x}] = \mathbb{E}^{\mathfrak{M};\operatorname{do}(X:=\mathbf{x})}[\operatorname{xor}(f(X), N_Y) \mid X = \mathbf{x}],$$

$$= \mathbb{E}^{\mathfrak{M}}[\operatorname{xor}(f(\mathbf{x}), N_Y)], \qquad (N_Y \perp \!\!\! \perp X^{\mathfrak{M};\operatorname{do}(X:=\mathbf{x})}.)$$

$$= \mathbb{E}^{\mathfrak{M}}[f(\mathbf{x}) + N_Y - 2f(\mathbf{x})N_Y], \qquad (Definition of xor.)$$

$$= f(\mathbf{x}) + \mathbb{E}^{\mathfrak{M}}[N_Y] - 2f(\mathbf{x})\mathbb{E}^{\mathfrak{M}}[N_Y],$$

$$= \left(1 - 2\mathbb{E}^{\mathfrak{M}}[N_Y]\right)f(\mathbf{x}) + \mathbb{E}^{\mathfrak{M}}[N_Y],$$

$$= 0.5f(\mathbf{x}) + 0.25. \qquad (N_Y \sim B(0.25).)$$

^{10.} Following the mechanisms $c_0 \to \mathbf{x}_1 \to \tilde{y}_2 \to y_3 \to c_4 \to \mathbf{x}_5$, we see that $(\mathbf{x}_4, y_4, c_4) = (\mathbf{x}_5, y_5, c_5)$ (same for $\tilde{y}_4 = \tilde{y}_5$).

Authors

Extended Abstract Track

- (a) Graph for generating colored-MNIST data.
- (b) Augmented graph exogenous variables explicitly shown.

Figure 7: A cyclic SEM perspective of the colored-MNIST data – an MNIST image N_X is assigned color C to produce a colored-MNIST image X. This is then passed through the ground-truth labeling function f to produce the true label \tilde{Y} . We flip this with probability 0.25 to produce the observed label Y, which in turn is flipped with probability e (at train time $e \in \{0.1, 0.2\}$ and e = 0.9 at test time) to produce the color C. These assignments are iteratively applied for any joint sample of the exogenous variables N_X, N_Y, N_C starting at arbitrary values of endogenous variables until convergence to the unique stationary point X, Y, C (and \tilde{Y}).

Because this is a binary classification task, we have

$$\operatorname{round}\!\left(\mathbb{E}^{\mathfrak{M};\operatorname{do}(X:=\mathbf{x})}[Y\,|\,X=\mathbf{x}]\right)=f(\mathbf{x}).$$

This is in contrast to the original DAG perspective of colored-MNIST shown in Fig. 6, where the connection to the estimation of the causal mechanism f is not immediately obvious. We argue that this is because the DAG in Fig. 6 is virtually equivalent to the reduced form of our structural form presented in Fig. 7.

Appendix H. Proofs

H.1. Proof of Proposition 1 – DA induced regularization

$$\Rightarrow \left\| \hat{\mathbf{h}}_{\mathrm{DA}_{G}+\mathrm{ERM}}^{\mathfrak{A}} \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2}$$

$$= \left\| \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[(GX)Y^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$= \left\| \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[(GX) \left(\mathbf{f}^{\top}X + \xi \right)^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$= \left\| \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[(GX) \left(\mathbf{f}^{\top}GX + \xi \right)^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$= \left\| \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[(GX)(GX)^{\top} \right] \mathbf{f} + \mathbb{E} \left[(GX)\xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$= \left\| \mathbf{f} + \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[(GX)\xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$= \left\| \mathbf{f} + \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[\left(X + \gamma \tilde{G} \right) \xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$= \left\| \mathbf{f} + \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$= \left\| \mathbf{f} + \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$= \left\| \mathbf{f} + \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$= \left\| \mathbf{f} + \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$= \left\| \mathbf{f} \|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2} + \left\| \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$= \left\| \mathbf{f} \|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2} + \left\| \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$= \left\| \mathbf{f} \|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2} + \left\| \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$= \left\| \mathbf{f} \|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2} + \left\| \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$= \left\| \mathbf{f} \|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2} + \left\| \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$= \left\| \mathbf{f} \|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2} + \left\| \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[(X \xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$= \left\| \mathbf{f} \|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2} + \left\| \mathbb{E} \left[(X \xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$\left\| \mathbf{f} \|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}^{\mathfrak{A}} + \mathbb{E} \left[(X \xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}}^{2},$$

$$\left\| \mathbf{f} \|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}^{\mathfrak{A}} + \mathbb{E} \left[(X \xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}^{\mathfrak{A}}^{\mathfrak{A}} + \mathbb{E} \left[(X \xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{A}}^{\mathfrak{A}} + \mathbb{E} \left[$$

where the first term does not depend on γ . The last term also does not depend on γ because

$$\Rightarrow \mathbf{f}^{\top} \mathbf{\Sigma}_{X}^{\mathfrak{A}} \mathbb{E} \left[\left(X + \gamma \tilde{G} \right) \left(X + \gamma \tilde{G} \right)^{\top} \right]^{-1} \mathbb{E} \left[X \xi^{\top} \right],$$

$$= \mathbf{f}^{\top} \mathbf{\Sigma}_{X}^{\mathfrak{A}} \mathbb{E} \left[\mathbf{\Sigma}_{X}^{\mathfrak{A}} + \gamma^{2} \mathbf{\Sigma}_{\tilde{G}}^{\mathfrak{A}} \right]^{-1} \mathbb{E} \left[X \xi^{\top} \right],$$

$$= \mathbf{f}^{\top} \mathbf{S}^{\top} \mathbf{S} \mathbb{E} \left[\mathbf{S}^{\top} \mathbf{S} + \gamma^{2} \mathbf{S}^{\top} \mathbf{D} \mathbf{S} \right]^{-1} \mathbb{E} \left[X \xi^{\top} \right], \qquad (\text{From Lemma 2.})$$

$$= \mathbf{f}^{\top} \mathbf{S}^{\top} \mathbf{S} \mathbf{S}^{-1} \mathbb{E} \left[\mathbf{I}_{m} + \gamma^{2} \mathbf{D} \right]^{-1} \mathbf{S}^{-\top} \mathbb{E} \left[X \xi^{\top} \right], \qquad (\mathbf{S}, \mathbf{S}^{\top} \text{ invertible.})$$

$$= \mathbf{f}^{\top} \mathbf{S}^{\top} \mathbb{E} \left[\mathbf{I}_{m} + \gamma^{2} \mathbf{D} \right]^{-1} \mathbf{S}^{-\top} \mathbb{E} \left[X \xi^{\top} \right],$$

$$= \mathbf{f}^{\top} \mathbf{S}^{\top} \mathbf{S}^{-\top} \mathbb{E} \left[X \xi^{\top} \right]. \qquad (\tilde{G} \in \text{null}(\mathbf{f}^{\top}) \Rightarrow \mathbf{f}^{\top} \mathbf{S}^{\top} \mathbf{D} = \mathbf{0}.)$$

Finally, for the middle term in Eq. (19) we can follow a similar approach as Theorem 1 to show that it is strictly decreasing in γ^2 , with equality iff

$$\mathbb{E}^{\mathfrak{A}}[GX|G] \perp \mathbb{E}^{\mathfrak{A}}[X|\xi]$$

H.2. Proof of Proposition 2 – IVL regression closed form solution in the linear case

The OLS solution for (X', Y') minimizes the following ERM risk

$$\Rightarrow \mathbb{E}\left[\left\|Y' - \mathbf{h}^{\top} X'\right\|^{2}\right]$$

$$= \mathbb{E}\left[\left\|aY + b\mathbb{E}[Y|Z] - \mathbf{h}^{\top}(aX + b\mathbb{E}[X|Z])\right\|^{2}\right], \quad \text{(Substitute in definitions of } X', Y'.\text{)}$$

$$= \mathbb{E}\left[\left\|a\left(Y - \mathbf{h}^{\top} X\right) + b\left(\mathbb{E}[Y|Z] - \mathbf{h}^{\top} \mathbb{E}[X|Z]\right)\right\|^{2}\right], \quad \text{(Distribute the subtraction.)}$$

$$= a^{2} \mathbb{E}\left[\left\|Y - \mathbf{h}^{\top} X\right\|^{2}\right] + b^{2} \mathbb{E}\left[\left\|\mathbb{E}[Y|Z] - \mathbf{h}^{\top} \mathbb{E}[X|Z]\right\|^{2}\right] \quad \text{(Expand squared norm.)}$$

$$+ 2ab \mathbb{E}\left[\left(Y - \mathbf{h}^{\top} X\right)^{\top} \left(\mathbb{E}[Y|Z] - \mathbf{h}^{\top} \mathbb{E}[X|Z]\right)\right]. \quad (20)$$

First we note that from the definitions of a, b we have

$$a^{2} = \sqrt{\alpha}, \qquad b^{2} + 2ab = \left(\sqrt{1+\alpha} - \sqrt{\alpha}\right)^{2} + 2\sqrt{\alpha}\left(\sqrt{1+\alpha} - \sqrt{\alpha}\right) = 1. \tag{21}$$

Now we evaluate the cross term in Eq. (20)

$$\Rightarrow \mathbb{E}\left[\left(Y - \mathbf{h}^{\top}X\right)^{\top}\left(\mathbb{E}[Y|Z] - \mathbf{h}^{\top}\mathbb{E}[X|Z]\right)\right]$$

$$= \mathbb{E}\left[\mathbb{E}\left[\left(Y - \mathbf{h}^{\top}X\right)^{\top}\left(\mathbb{E}[Y|Z] - \mathbf{h}^{\top}\mathbb{E}[X|Z]\right)\Big|Z\right]\right], \quad \text{(Law of iterated expectation.)}$$

$$= \mathbb{E}\left[\mathbb{E}\left[\left(Y - \mathbf{h}^{\top}X\right)^{\top}\Big|Z\right]\left(\mathbb{E}[Y|Z] - \mathbf{h}^{\top}\mathbb{E}[X|Z]\right)\right] \quad \text{(Taking out what is known; Eq. (15).)}$$

$$= \mathbb{E}\left[\left(\mathbb{E}[Y|Z] - \mathbf{h}^{\top}\mathbb{E}[X|Z]\right)^{\top}\left(\mathbb{E}[Y|Z] - \mathbf{h}^{\top}\mathbb{E}[X|Z]\right)\right]$$

$$= \mathbb{E}\left[\left\|\mathbb{E}[Y|Z] - \mathbf{h}^{\top}\mathbb{E}[X|Z]\right\|^{2}\right].$$

Substituting this back in Eq. (20) we get

$$\Rightarrow \mathbb{E}\left[\left\|Y' - \mathbf{h}^{\top} X'\right\|^{2}\right]$$

$$= a^{2} \mathbb{E}\left[\left\|Y - \mathbf{h}^{\top} X\right\|^{2}\right] + \left(b^{2} + 2ab\right) \mathbb{E}\left[\left\|\mathbb{E}[Y|Z] - \mathbf{h}^{\top} \mathbb{E}[X|Z]\right\|^{2}\right],$$

$$= \alpha \mathbb{E}\left[\left\|Y - \mathbf{h}^{\top} X\right\|^{2}\right] + \mathbb{E}\left[\left\|\mathbb{E}[Y|Z] - \mathbf{h}^{\top} \mathbb{E}[X|Z]\right\|^{2}\right], \qquad \text{(From Eq. (21).)}$$

$$= \alpha R_{\text{ERM}}^{\mathfrak{M}}(\mathbf{h}) + R_{\text{IV}}^{\mathfrak{M}}(\mathbf{h}) - \mathbb{E}[\mathbb{V}[Y|Z]], \qquad \text{(From Eq. (13).)}$$

$$= R_{\text{IVI}}^{\mathfrak{M}}\left(\mathbf{h}\right) - \mathbb{E}[\mathbb{V}[Y|Z]].$$

H.3. Proof of Proposition 3 – Existence of interventional distribution far a DA

Proposition 3 (unique stationary interventional distribution) In SEM \mathfrak{A} from Eq. (7), given any $(\mathbf{g}, \mathbf{c}, \mathbf{n}_X, \mathbf{n}_Y) \sim P_{G,C,N_X,N_Y}^{\mathfrak{A}}$, if for all $(\mathbf{x}_0, \mathbf{y}_0) \in \mathcal{X} \times \mathcal{Y}$ the unique limits

$$\mathbf{x}^{\mathfrak{A}} \coloneqq \lim_{t \to \infty} \mathbf{x}_{t}^{\mathfrak{A}} = \lim_{t \to \infty} \tau \left(\mathbf{y}_{t-1}^{\mathfrak{A}}, \mathbf{c}, \mathbf{n}_{X} \right),$$

$$\mathbf{y}^{\mathfrak{A}} \coloneqq \lim_{t \to \infty} \mathbf{y}_{t}^{\mathfrak{A}} = \lim_{t \to \infty} f \left(\mathbf{x}_{t-1}^{\mathfrak{A}} \right) + \epsilon(\mathbf{c}) + \mathbf{n}_{Y}$$

exist, then in \mathfrak{A} ; do($\tau := \mathbf{g}\tau$) the unique limits

$$\begin{split} \mathbf{x}^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)} &:= \lim_{t \to \infty} \mathbf{x}_t^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)} = \lim_{t \to \infty} \mathbf{g}\tau \Big(\mathbf{y}_{t-1}^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)}, \mathbf{c}, \mathbf{n}_X \Big) = \mathbf{g}\mathbf{x}^{\mathfrak{A}}, \\ \mathbf{y}^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)} &:= \lim_{t \to \infty} \mathbf{y}_t^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)} = \lim_{t \to \infty} f\Big(\mathbf{x}_{t-1}^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)} \Big) + \epsilon(\mathbf{c}) + \mathbf{n}_Y = \mathbf{y}^{\mathfrak{A}} \end{split}$$

also exist.

Proof First we try to show that

$$\mathbf{y}_t^{\mathfrak{A};\operatorname{do}(\tau := \mathbf{g}\tau)} = \mathbf{y}_t^{\mathfrak{A}}.$$
 (22)

For the base case, we have by construction

$$\mathbf{y}_0^{\mathfrak{A};\operatorname{do}(\tau \coloneqq \mathbf{g}\tau)} \coloneqq \mathbf{y}_0 \eqqcolon \mathbf{y}_0^{\mathfrak{A}}.$$

For the step case, assuming that $\mathbf{y}_t^{\mathfrak{A}; \text{do}(\tau := \mathbf{g}\tau)} = \mathbf{y}_t^{\mathfrak{A}}$, we have 11,

$$\mathbf{y}_{t+2}^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)} = f\left(\mathbf{x}_{t+1}^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)}\right) + \epsilon(\mathbf{c}) + \mathbf{n}_{Y},$$

$$= f\left(\mathbf{g}\tau\left(\mathbf{y}_{t}^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)}, \mathbf{c}, \mathbf{n}_{X}\right)\right) + \epsilon(\mathbf{c}) + \mathbf{n}_{Y},$$

$$= f\left(\tau\left(\mathbf{y}_{t}^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)}, \mathbf{c}, \mathbf{n}_{X}\right)\right) + \epsilon(\mathbf{c}) + \mathbf{n}_{Y}, \qquad \text{(Invariance of } f \text{ to } \mathbf{g}.\text{)}$$

$$= f\left(\tau\left(\mathbf{y}_{t}^{\mathfrak{A}}, \mathbf{c}, \mathbf{n}_{X}\right)\right) + \epsilon(\mathbf{c}) + \mathbf{n}_{Y}, \qquad \text{(Assumption } \mathbf{y}_{t}^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)} = \mathbf{y}_{t}^{\mathfrak{A}}.\text{)}$$

$$= f\left(\mathbf{x}_{t+1}^{\mathfrak{A}}\right) + \epsilon(\mathbf{c}) + \mathbf{n}_{Y},$$

$$= \mathbf{y}_{t+1}^{\mathfrak{A}} \circ \cdot$$

Hence, we have shown that Eq. (22) holds for all even t. For odd t, we simply replace t = 0 with t = 1 in the base case

$$\mathbf{y}_{1}^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)} = f\left(\mathbf{x}_{0}^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)}\right) + \epsilon(\mathbf{c}) + \mathbf{n}_{Y},$$

$$= f\left(\mathbf{x}_{0}^{\mathfrak{A}}\right) + \epsilon(\mathbf{c}) + \mathbf{n}_{Y}, \qquad \text{(Definitions } \mathbf{x}_{0}^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)} := \mathbf{x}_{0} =: \mathbf{x}_{0}^{\mathfrak{A}}.)$$

$$= \mathbf{y}_{1}^{\mathfrak{A}},$$

^{11.} Note that here the step size for proof by induction would be $\Delta t = 2$ since \mathbf{y}_t precedes \mathbf{y}_{t+2} . Similar is the case for \mathbf{x}_t as well.

Authors

Extended Abstract Track

We have now finally shown that Eq. (22) holds for all $t \geq 0$.

Next, it is now relatively straightforward to show that for any t > 0, we have

$$\mathbf{x}_{t}^{\mathfrak{A};\text{do}(\tau:=\mathbf{g}\tau)} = \mathbf{g}\tau\left(\mathbf{y}_{t-1}^{\mathfrak{A};\text{do}(\tau:=\mathbf{g}\tau)}, \mathbf{c}, \mathbf{n}_{X}\right),$$

$$= \mathbf{g}\tau\left(\mathbf{y}_{t-1}^{\mathfrak{A}}, \mathbf{c}, \mathbf{n}_{X}\right), \qquad \text{(Follows from Eq. (22).)}$$

$$= \mathbf{g}\mathbf{x}_{t}^{\mathfrak{A}}. \qquad (23)$$

Finally, by applying limit as $t \to \infty$ to both sides of Eq. (22) and Eq. (23), we get

$$\mathbf{y}^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)} = \lim_{t \to \infty} \mathbf{y}_{t}^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)} = \lim_{t \to \infty} \mathbf{y}_{t}^{\mathfrak{A}} = \mathbf{y}^{\mathfrak{A}},$$

$$\mathbf{x}^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)} = \lim_{t \to \infty} \mathbf{x}_{t}^{\mathfrak{A};\operatorname{do}(\tau:=\mathbf{g}\tau)} = \lim_{t \to \infty} \mathbf{g}\mathbf{x}_{t}^{\mathfrak{A}} = \mathbf{g}\lim_{t \to \infty} \mathbf{x}_{t}^{\mathfrak{A}} = \mathbf{g}\mathbf{x}^{\mathfrak{A}},$$
(24)

where the limit can be moved past \mathbf{g} in Eq. (24) because \mathbf{g} is assumed continuous in its domain.

H.4. Proof of Theorem 1 – Causal estimation with DA+ERM

$$\begin{split} &\Rightarrow \left\| \hat{\mathbf{h}}_{\mathrm{D}A_{G}^{+}\mathrm{ERM}}^{\mathrm{R}} - \mathbf{f} \right\|_{\boldsymbol{\Sigma}_{X}^{\mathrm{R}}} = \left\| \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[(GX)Y^{\top} \right] - \mathbf{f} \right\|_{\boldsymbol{\Sigma}_{X}^{\mathrm{R}}}, \\ &= \left\| \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[(GX) \left(\mathbf{f}^{\top}X + \xi \right)^{\top} \right] - \mathbf{f} \right\|_{\boldsymbol{\Sigma}_{X}^{\mathrm{R}}}, & \text{(Structural eq. of } Y.) \\ &= \left\| \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[(GX) \left(\mathbf{f}^{\top}(GX) + \xi \right)^{\top} \right] - \mathbf{f} \right\|_{\boldsymbol{\Sigma}_{X}^{\mathrm{R}}}, & \text{(Using } \mathcal{G}\text{-invariance of } \mathbf{f}.) \\ &= \left\| \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[(GX)\xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathrm{R}}}, & \\ &= \left\| \mathbb{E} \left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E} \left[(GX)\xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathrm{R}}}, & \text{(\bar{G} := } \mathbb{E}[GX|G] = $\gamma \cdot \Gamma^{\top}G.$) \\ &= \left\| \left(\mathbb{E} \left[XX^{\top} \right] + \mathbb{E} \left[\bar{G}\bar{G}^{\top} \right] \right)^{-1} \mathbb{E} \left[X\xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathrm{R}}}, & \text{(Using } \bar{\mathcal{G}} \perp X, \xi.) \\ &= \left\| \left(\mathbb{E} \left[XX^{\top} \right] + \mathbb{E} \left[\bar{G}\bar{G}^{\top} \right] \right)^{-1} \mathbb{E} \left[X\xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathrm{R}}}, & \text{(Using } \bar{\mathcal{G}} \perp X, \xi.) \\ &= \left\| \left(\mathbb{E} \left[XX^{\top} \right] + \mathbb{E} \left[\bar{G}\bar{G}^{\top} \right] \right)^{-1} \mathbb{E} \left[X\xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathrm{R}}}, & \text{(Using } \bar{\mathcal{G}} \perp X, \xi.) \\ &= \left\| \left(\mathbb{E} \left[XX^{\top} \right] - \mathbb{E} \left[X\xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathrm{R}}}, & \text{(Switch to } \ell_{2} \text{ norm.} \right. \\ &= \left\| \mathbb{E} \left[XX^{\top} \right]^{-1} \mathbb{E} \left[X\xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathrm{R}}}, & \text{(Substitute in } \mathbf{I}_{m} = \mathbf{S}\mathbf{S}^{-1}. \\ &= \left\| \mathbb{E} \left[XX^{\top} \right]^{-1} \mathbb{E} \left[X\xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathrm{R}}}, & \text{(Substitute in } \mathbf{\Sigma}_{M}^{\mathrm{R}} := \mathbb{E}^{2M} \left[XX^{\top} \right] = \mathbf{S}^{\mathrm{T}}\mathbf{S}. \\ &= \left\| \mathbb{E} \left[XX^{\top} \right]^{-1} \mathbb{E} \left[X\xi^{\top} \right] + \mathbb{E} \left[X\xi^{\top} \right] \right\|_{\boldsymbol{\Sigma}_{X}^{\mathrm{R}}}, & \text{(Add and subtract } \mathbf{f}. \\ &= \left\| \mathbb{E} \left[XX^{\top} \right]^{-1} \mathbb{E} \left[X\left(\mathbf{f}^{\top}X + \xi \right)^{\top} \right] - \mathbf{f} \right\|_{\boldsymbol{\Sigma}_{X}^{\mathrm{R}}}, & \text{(Linearity of expectation.)} \\ &= \left\| \mathbb{E} \left[XX^{\top} \right]^{-1} \mathbb{E} \left[XY^{\top} \right] - \mathbf{f} \right\|_{\boldsymbol{\Sigma}^{\mathrm{R}}} & \mathbb{E}^{\mathrm{R}} - \mathbf{f} \right\|_{\boldsymbol{\Sigma}_{X}^{\mathrm{R}}}, & \text{(Structural eq. of } Y. \\ \end{pmatrix}$$

where inequality Eq. (25) holds because **D** is non-negative diagonal. Also, inequality Eq. (25) only holds with equality iff $\mathbf{S}^{-\top}\mathbb{E}\left[X\xi^{\top}\right]$ is in kernel of **D**. Or equivalently, iff $\mathbb{E}\left[X\xi^{\top}\right]$ is in the kernel of $\mathbf{S}^{\top}\mathbf{D}\mathbf{S} = \mathbf{\Sigma}_{\tilde{G}}$, which from Lemma 1 holds iff $\mathbb{E}^{\mathfrak{M}}[GX|G] \perp \mathbb{E}^{\mathfrak{M}}[X|\xi]$ a.s.

H.5. Proof of Theorem 2 – Robust prediction with IVL regression

Write X in terms of the exogenous variables C, Z, N_X, N_Y using the reduced form from Lemma 3 as

$$X = \tilde{Z} + \tilde{C} + \tilde{N},\tag{26}$$

where for readability we represent

$$\tilde{Z} \coloneqq \mathbf{M}_{m \times m} \mathbf{\Gamma}^{\top} Z, \qquad \qquad \tilde{C} \coloneqq \mathbf{M} \begin{bmatrix} \mathbf{T}^{\top} \\ \boldsymbol{\epsilon}^{\top} \end{bmatrix} C, \qquad \qquad \tilde{N} \coloneqq \sigma \cdot \mathbf{M} \begin{bmatrix} N_X \\ N_Y \end{bmatrix},$$

with

$$\mathbf{M} \coloneqq egin{bmatrix} \mathbf{M}_{m imes m} & \mathbf{M}_{m imes 1} \ \mathbf{M}_{1 imes m} & \mathbf{M}_{1 imes 1} \end{bmatrix} = egin{bmatrix} \mathbf{I}_m & -oldsymbol{ au}^ op \ -\mathbf{f}^ op & 1 \end{bmatrix}^{-1}.$$

Now, we start by writing the ERM objective under the intervention do $(\Gamma^{\top}(\cdot) \coloneqq \zeta)$ as

$$\Rightarrow R_{\text{ERM}}^{\mathfrak{M};\text{do}\left(\Gamma^{\top}(\cdot):=\zeta\right)}(\mathbf{h})$$

$$= \mathbb{E}^{\mathfrak{M};\text{do}\left(\Gamma^{\top}(\cdot):=\zeta\right)} \left[\left\| Y - \mathbf{h}^{\top} X \right\|^{2} \right],$$

$$= \mathbb{E}^{\mathfrak{M};\text{do}\left(\Gamma^{\top}(\cdot):=\zeta\right)} \left[\left\| \xi + (\mathbf{f} - \mathbf{h})^{\top} \left(\tilde{Z} + \tilde{C} + \tilde{N} \right) \right\|^{2} \right], \qquad (Y \text{ structural form } \& \text{ Eq. (26).})$$

$$= \mathbb{E}^{\mathfrak{M};\text{do}\left(\Gamma^{\top}(\cdot):=\zeta\right)} \left[\left\| \xi + (\mathbf{f} - \mathbf{h})^{\top} \left(\mathbf{M}_{m \times m} \zeta + \tilde{C} + \tilde{N} \right) \right\|^{2} \right], \quad (\tilde{Z} \& \text{ intervention definition.})$$

$$= \mathbb{E}^{\mathfrak{M};\text{do}\left(\Gamma^{\top}(\cdot):=\zeta\right)} \left[\left\| \xi + (\mathbf{f} - \mathbf{h})^{\top} \left(\tilde{C} + \tilde{N} \right) + (\mathbf{f} - \mathbf{h})^{\top} \mathbf{M}_{m \times m} \zeta \right\|^{2} \right], \quad (\text{Define } \mathbf{h}'^{\top} := (\mathbf{f} - \mathbf{h})^{\top} \mathbf{M}_{m \times m}.)$$

$$= \mathbb{E}^{\mathfrak{M};\text{do}\left(\Gamma^{\top}(\cdot):=\zeta\right)} \left[\left\| \xi + (\mathbf{f} - \mathbf{h})^{\top} \left(\tilde{C} + \tilde{N} \right) + \mathbf{h}'^{\top} \zeta \right\|^{2} \right], \quad (\text{Pollows from exogeneity of } \zeta \text{ under intervention, } \Rightarrow \text{ cross term zeros-out.})$$

$$= \mathbb{E}^{\mathfrak{M};\text{do}\left(\Gamma^{\top}(\cdot):=\mathbf{0}_{m}\right)} \left[\left\| Y - \mathbf{h}^{\top} X \right\|^{2} \right] + \mathbb{E}^{\mathfrak{M};\text{do}\left(\Gamma^{\top}(\cdot):=\zeta\right)} \left[\left\| \mathbf{h}'^{\top} \zeta \right\|^{2} \right], \quad (27)$$

$$= \mathbb{E}^{\mathfrak{M};\text{do}\left(\Gamma^{\top}(\cdot):=\mathbf{0}_{m}\right)} \left[\left\| Y - \mathbf{h}^{\top} X \right\|^{2} \right] + \left\| \mathbf{h}'^{\top} \zeta \right\|^{2},$$

$$= \mathbb{E}^{\mathfrak{M};\text{do}\left(\Gamma^{\top}(\cdot):=\mathbf{0}_{m}\right)} \left[\left\| Y - \mathbf{h}^{\top} X \right\|^{2} \right] + \text{tr}\left(\zeta^{\top} \mathbf{h}' \mathbf{h}'^{\top} \zeta\right),$$

$$= \mathbb{E}^{\mathfrak{M};\text{do}\left(\Gamma^{\top}(\cdot):=\mathbf{0}_{m}\right)} \left[\left\| Y - \mathbf{h}^{\top} X \right\|^{2} \right] + \text{tr}\left(\zeta^{\top} \mathbf{h}' \mathbf{h}'^{\top} \zeta\right),$$

$$= \mathbb{E}^{\mathfrak{M};\text{do}\left(\Gamma^{\top}(\cdot):=\mathbf{0}_{m}\right)} \left[\left\| Y - \mathbf{h}^{\top} X \right\|^{2} \right] + \text{tr}\left(\zeta^{\top} \mathbf{h}' \mathbf{h}'^{\top} \zeta\right),$$

Now, note that the maximum of the trace term over $\zeta \in \mathcal{P}_{\alpha}$ gives

$$\Rightarrow \max_{\boldsymbol{\zeta} \in \mathcal{P}_{\alpha}} \ \mathrm{tr} \Big(\mathbf{h}'^{\top} \boldsymbol{\zeta} \boldsymbol{\zeta}^{\top} \mathbf{h}' \Big),$$

$$= \left(\frac{1}{\alpha} + 1\right) \operatorname{tr}\left(\mathbf{h}'^{\top}\left(\mathbf{\Gamma}^{\top}\mathbb{E}^{\mathfrak{M}}\left[ZZ^{\top}\right]\mathbf{\Gamma}\right)\mathbf{h}'\right), \quad \text{(Linearity of trace and definition of } \mathcal{P}_{\alpha}.)$$

$$= \left(\frac{1}{\alpha} + 1\right) \mathbb{E}^{\mathfrak{M}}\left[\operatorname{tr}\left(\mathbf{h}'^{\top}\mathbf{\Gamma}^{\top}ZZ^{\top}\mathbf{\Gamma}\mathbf{h}'\right)\right], \quad \text{(Linearity of expectation.)}$$

$$= \left(\frac{1}{\alpha} + 1\right) \mathbb{E}^{\mathfrak{M}}\left[\operatorname{tr}\left(Z^{\top}\mathbf{\Gamma}\mathbf{h}'\mathbf{h}'^{\top}\mathbf{\Gamma}^{\top}Z\right)\right], \quad \text{(Cyclic property of trace.)}$$

$$= \left(\frac{1}{\alpha} + 1\right) \mathbb{E}^{\mathfrak{M}}\left[\left\|\mathbf{h}'^{\top}\mathbf{\Gamma}^{\top}Z\right\|^{2}\right], \quad \text{(Substitute in definition of } \mathbf{h}'^{\top}.)$$

$$= \left(\frac{1}{\alpha} + 1\right) \mathbb{E}^{\mathfrak{M}}\left[\left\|(\mathbf{f} - \mathbf{h})^{\top}\mathbf{M}_{m \times m}\mathbf{\Gamma}^{\top}Z\right\|^{2}\right]. \quad \text{(Definition of } \tilde{Z}.)$$

We can now substitute this in while maximizing both sides of Eq. (28) over interventions $\zeta \in \mathcal{P}_{\alpha}$ as

$$\begin{split} &\Rightarrow \max_{\boldsymbol{\zeta} \in \mathcal{P}_{\alpha}} R_{\mathrm{ERM}}^{\mathfrak{M};\mathrm{do}\left(\boldsymbol{\Gamma}^{\top}(\cdot) \coloneqq \mathbf{0}_{m}\right)}(\mathbf{h}) \\ &= \mathbb{E}^{\mathfrak{M};\mathrm{do}\left(\boldsymbol{\Gamma}^{\top}(\cdot) \coloneqq \mathbf{0}_{m}\right)} \left[\left\| \boldsymbol{Y} - \mathbf{h}^{\top} \boldsymbol{X} \right\|^{2} \right] + \max_{\boldsymbol{\zeta} \in \mathcal{P}_{\alpha}} \mathrm{tr}\left(\mathbf{h}'^{\top} \boldsymbol{\zeta} \boldsymbol{\zeta}^{\top} \mathbf{h}'\right), \text{ (First term does not have } \boldsymbol{\zeta}.) \\ &= \mathbb{E}^{\mathfrak{M};\mathrm{do}\left(\boldsymbol{\Gamma}^{\top}(\cdot) \coloneqq \mathbf{0}_{m}\right)} \left[\left\| \boldsymbol{Y} - \mathbf{h}^{\top} \boldsymbol{X} \right\|^{2} \right] + \left(\frac{1}{\alpha} + 1 \right) \mathbb{E}^{\mathfrak{M}} \left[\left\| (\mathbf{f} - \mathbf{h})^{\top} \tilde{\boldsymbol{Z}} \right\|^{2} \right], \\ &= \mathbb{E}^{\mathfrak{M}} \left[\left\| \boldsymbol{Y} - \mathbf{h}^{\top} \boldsymbol{X} \right\|^{2} \right] + \frac{1}{\alpha} \mathbb{E}^{\mathfrak{M}} \left[\left\| (\mathbf{f} - \mathbf{h})^{\top} \mathbb{E}[\boldsymbol{X} \mid \boldsymbol{Z}] \right\|^{2} \right], \\ &= \mathbb{E}^{\mathfrak{M}} \left[\left\| \boldsymbol{Y} - \mathbf{h}^{\top} \boldsymbol{X} \right\|^{2} \right] + \frac{1}{\alpha} \mathbb{E}^{\mathfrak{M}} \left[\left\| (\mathbf{f} - \mathbf{h})^{\top} \mathbb{E}[\boldsymbol{X} \mid \boldsymbol{Z}] \right\|^{2} \right], \\ &= \mathbb{E}^{\mathfrak{M}} \left[\left\| \boldsymbol{Y} - \mathbf{h}^{\top} \boldsymbol{X} \right\|^{2} \right] + \frac{1}{\alpha} \mathbb{E}^{\mathfrak{M}} \left[\left\| \mathbb{E} \left[\mathbf{f}^{\top} \boldsymbol{X} \mid \boldsymbol{Z} \right] - \mathbf{h}^{\top} \mathbb{E}[\boldsymbol{X} \mid \boldsymbol{Z}] \right\|^{2} \right], \\ &= \mathbb{E}^{\mathfrak{M}} \left[\left\| \boldsymbol{Y} - \mathbf{h}^{\top} \boldsymbol{X} \right\|^{2} \right] + \frac{1}{\alpha} \mathbb{E}^{\mathfrak{M}} \left[\left\| \mathbb{E}[\boldsymbol{Y} \mid \boldsymbol{Z}] - \mathbf{h}^{\top} \mathbb{E}[\boldsymbol{X} \mid \boldsymbol{Z}] \right\|^{2} \right], \\ &= \mathbb{E}^{\mathfrak{M}} \left[\left\| \boldsymbol{Y} - \mathbf{h}^{\top} \boldsymbol{X} \right\|^{2} \right] + \frac{1}{\alpha} \mathbb{E}^{\mathfrak{M}} \left[\left\| \mathbb{E}[\boldsymbol{Y} \mid \boldsymbol{Z}] - \mathbf{h}^{\top} \mathbb{E}[\boldsymbol{X} \mid \boldsymbol{Z}] \right\|^{2} \right], \\ &= \mathbb{E}^{\mathfrak{M}} \left[\left\| \boldsymbol{Y} - \mathbf{h}^{\top} \boldsymbol{X} \right\|^{2} \right] + \mathbb{E}^{\mathfrak{M}} \left[\left\| \mathbb{E}[\boldsymbol{Y} \mid \boldsymbol{Z}] - \mathbf{h}^{\top} \mathbb{E}[\boldsymbol{X} \mid \boldsymbol{Z}] \right\|^{2} \right], \\ &= \mathbb{E}^{\mathfrak{M}} \left[\left\| \boldsymbol{Y} - \mathbf{h}^{\top} \boldsymbol{X} \right\|^{2} \right] + \mathbb{E}^{\mathfrak{M}} \left[\left\| \mathbb{E}[\boldsymbol{Y} \mid \boldsymbol{Z}] - \mathbf{h}^{\top} \mathbb{E}[\boldsymbol{X} \mid \boldsymbol{Z}] \right\|^{2} \right], \\ &= \mathbb{E}^{\mathfrak{M}} \left[\left\| \boldsymbol{Y} - \mathbf{h}^{\top} \boldsymbol{X} \right\|^{2} \right] + \mathbb{E}^{\mathfrak{M}} \left[\left\| \mathbb{E}[\boldsymbol{Y} \mid \boldsymbol{Z}] - \mathbf{h}^{\top} \mathbb{E}[\boldsymbol{X} \mid \boldsymbol{Z}] \right\|^{2} \right], \\ &= \mathbb{E}^{\mathfrak{M}} \left[\left\| \boldsymbol{Y} - \mathbf{h}^{\top} \boldsymbol{X} \right\|^{2} \right] + \mathbb{E}^{\mathfrak{M}} \left[\mathbb{E}[\boldsymbol{Y} \mid \boldsymbol{Z}] - \mathbf{h}^{\top} \mathbb{E}[\boldsymbol{X} \mid \boldsymbol{Z}] \right]^{2} \right], \end{aligned} \quad \text{(Inverse step of Eq. (27).)} \\ &= \mathbb{E}^{\mathfrak{M}} \left[\left\| \boldsymbol{Y} - \mathbf{h}^{\top} \boldsymbol{X} \right\|^{2} \right] + \mathbb{E}^{\mathfrak{M}} \left[\mathbb{E}[\boldsymbol{Y} \mid \boldsymbol{Z}] - \mathbf{h}^{\top} \mathbb{E}[\boldsymbol{X} \mid \boldsymbol{Z}] \right]^{2} \right], \end{aligned} \quad \text{(Inverse step of Eq. (27).)} \\ &= \mathbb{E}^{\mathfrak{M}} \left[\mathbb{E}[\boldsymbol{Y} \mid \boldsymbol{Z}] \right] + \mathbb{E}^{\mathfrak{M}} \left[\mathbb{E}[\boldsymbol{Z} \mid \boldsymbol{Z}] \right] + \mathbb{E}^{\mathfrak{M}} \left[\mathbb{E}[\boldsymbol{Z} \mid \boldsymbol{Z}] \right] + \mathbb{E}^{\mathfrak{M}} \left[\mathbb{E}[\boldsymbol{Z} \mid \boldsymbol{Z}] \right] \right]$$

H.6. Proof of Theorem 3 – Causal estimation with IVL regression

For $\hat{\mathbf{h}}_{\mathrm{IVL}_{\alpha}}^{\mathfrak{M}}$, we have from Proposition 2

$$\left\|\hat{\mathbf{h}}_{\mathrm{IVL}_{\alpha}}^{\mathfrak{M}} - \mathbf{f}\right\|_{\boldsymbol{\Sigma}_{X}^{\mathfrak{M}}}^{2} = \left\|\mathbb{E}\left[X'{X'}^{\top}\right]^{-1}\mathbb{E}\left[X'{Y'}^{\top}\right] - \mathbf{f}\right\|_{\boldsymbol{\Sigma}_{Y}^{\mathfrak{M}}}^{2}.$$

Note that we have

$$\Rightarrow \mathbb{E}\left[X'Y'^{\top}\right]$$

$$= \mathbb{E}\left[X'(aY + b\mathbb{E}[Y|Z])^{\top}\right],$$

$$= \mathbb{E}\left[X'\left(aY + b\mathbf{f}^{\top}\mathbb{E}[X|Z]\right)^{\top}\right],$$

$$= \mathbb{E}\left[X'\left(aY + b\mathbf{f}^{\top}\mathbb{E}[X|Z]\right)^{\top}\right],$$

$$= \mathbb{E}\left[X'\left(a\mathbf{f}^{\top}X + a\xi + b\mathbf{f}^{\top}\mathbb{E}[X|Z]\right)^{\top}\right],$$

$$= \mathbb{E}\left[X'\left(\mathbf{f}^{\top}X' + a\xi\right)^{\top}\right],$$

$$= \mathbb{E}\left[X'\mathbf{f}^{\top}X' + a\xi\right],$$

$$= \mathbb{E}\left[X'X'^{\top}\mathbf{f} + aX'\xi^{\top}\right],$$

$$= \mathbb{E}\left[X'X'^{\top}\mathbf{f} + a\mathbb{E}\left[X'\xi^{\top}\right],$$

$$= \mathbb{E}\left[X'X'^{\top}\mathbf{f} + a\mathbb{E}\left[X'\xi^{\top}\right],$$

$$= \mathbb{E}\left[X'X'^{\top}\mathbf{f} + a\mathbb{E}\left[X\xi^{\top}\right],$$

$$= \mathbb{E}\left[X'^{\top}\mathbf{f} + a\mathbb{E}\left[X\xi^{\top}\right],$$

$$= \mathbb{E}\left[X'^{\top}\mathbf{f} + a\mathbb{E}\left[X\xi^{\top}\right],$$

$$= \mathbb{E$$

We also see that

$$\Rightarrow \mathbb{E}\left[X'X'^{\top}\right]$$

$$= \mathbb{E}\left[(aX + b\mathbb{E}[X|Z])(aX + b\mathbb{E}[X|Z])^{\top}\right],$$

$$= \mathbb{E}\left[\left(aX + b\tilde{Z}\right)\left(aX + b\tilde{Z}\right)^{\top}\right], \qquad (\text{Set } \tilde{Z} := \mathbb{E}[X|Z] \text{ for brevity.})$$

$$= a^{2}\mathbb{E}\left[XX^{\top}\right] + b^{2}\mathbb{E}\left[\tilde{Z}\tilde{Z}^{\top}\right] + ab\mathbb{E}\left[X\tilde{Z}^{\top}\right] + ab\mathbb{E}\left[\tilde{Z}X^{\top}\right],$$

$$= a^{2}\mathbb{E}\left[XX^{\top}\right] + (b^{2} + 2ab)\boldsymbol{\Sigma}_{\tilde{Z}}, \qquad (\text{Because } \mathbb{E}\left[X\tilde{Z}^{\top}\right] = \boldsymbol{\Sigma}_{\tilde{Z}}.)$$

$$= \alpha\mathbb{E}\left[XX^{\top}\right] + \boldsymbol{\Sigma}_{\tilde{Z}}, \qquad (30)$$

where we substituted in Eq. (21) in Eq. (30).

Finally, we now have

$$\Rightarrow \left\|\hat{\mathbf{h}}_{\mathrm{IVL}_{\alpha}}^{\mathfrak{M}} - \mathbf{f}\right\|_{\boldsymbol{\Sigma}_{\mathbf{v}}^{\mathfrak{M}}}^{2}$$

$$= \left\| \mathbb{E} \left[X'X'^{\top} \right]^{-1} \mathbb{E} \left[X'Y'^{\top} \right] - \mathbf{f} \right\|_{\Sigma_{X}^{m}}^{2},$$

$$= \left\| \mathbb{E} \left[X'X'^{\top} \right]^{-1} \left(\mathbb{E} \left[X'X'^{\top} \right] \mathbf{f} + \alpha \mathbb{E} \left[X \xi^{\top} \right] \right) - \mathbf{f} \right\|_{\Sigma_{X}^{m}}^{2},$$

$$= \left\| \mathbf{f} + \alpha \mathbb{E} \left[X'X'^{\top} \right]^{-1} \mathbb{E} \left[X \xi^{\top} \right] - \mathbf{f} \right\|_{\Sigma_{X}^{m}}^{2},$$

$$= \left\| \alpha \mathbb{E} \left[X'X'^{\top} \right]^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{\Sigma_{X}^{m}}^{2},$$

$$= \left\| \alpha \left(\alpha \mathbb{E} \left[XX^{\top} \right] + \mathbf{\Sigma}_{\tilde{Z}} \right)^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{\Sigma_{X}^{m}}^{2},$$

$$= \left\| \left(\mathbf{S}^{\top} \mathbf{S} + \frac{1}{\alpha} \mathbf{S}^{\top} \mathbf{D} \mathbf{S} \right)^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{S^{\top} \mathbf{S}}^{2},$$

$$= \left\| \left(\mathbf{I}_{m} + \frac{1}{\alpha} \mathbf{D} \right)^{-1} \mathbf{S}^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{S^{\top} \mathbf{S}}^{2},$$

$$= \left\| \left(\mathbf{I}_{m} + \frac{1}{\alpha} \mathbf{D} \right)^{-1} \mathbf{S}^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{S^{\top} \mathbf{S}}^{2},$$

$$= \left\| \left(\mathbf{I}_{m} + \frac{1}{\alpha} \mathbf{D} \right)^{-1} \mathbf{S}^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{S^{\top} \mathbf{S}}^{2},$$

$$= \left\| \left(\mathbf{I}_{m} + \frac{1}{\alpha} \mathbf{D} \right)^{-1} \mathbf{S}^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{S^{\top} \mathbf{S}}^{2},$$

$$= \left\| \mathbf{S}^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{S^{-1} \mathbf{S}}^{2},$$

$$= \left\| \mathbf{S}^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{S^{-1} \mathbf{S}}^{2},$$

$$= \left\| \mathbf{S}^{-1} \mathbf{S}^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{S^{-1} \mathbf{S}}^{2},$$

$$= \left\| \mathbf{S}^{-1} \mathbf{S}^{-1} \mathbb{E} \left[X \xi^{\top} \right] \right\|_{S^{-1} \mathbf{S}}^{2},$$

$$= \left\| \mathbf{E} \left[X X^{\top} \right]^{-1} \mathbb{E} \left[X \xi^{\top} \right] - \mathbf{f} \right\|_{\Sigma_{\infty}^{m}}^{2},$$

$$= \left\| \mathbf{E} \left[X X^{\top} \right]^{-1} \left(\mathbb{E} \left[X X^{\top} \right] \mathbf{f} + \mathbb{E} \left[X \xi^{\top} \right] \right) - \mathbf{f} \right\|_{\Sigma_{\infty}^{m}}^{2},$$

$$= \left\| \mathbb{E} \left[X X^{\top} \right]^{-1} \mathbb{E} \left[X \left(\mathbf{f}^{\top} X + \xi \right)^{\top} \right] - \mathbf{f} \right\|_{\Sigma_{\infty}^{m}}^{2},$$

$$= \left\| \mathbb{E} \left[X X^{\top} \right]^{-1} \mathbb{E} \left[X \left(\mathbf{f}^{\top} X + \xi \right)^{\top} \right] - \mathbf{f} \right\|_{\Sigma_{\infty}^{m}}^{2},$$

$$= \left\| \mathbb{E} \left[X X^{\top} \right]^{-1} \mathbb{E} \left[X \left(\mathbf{f}^{\top} X + \xi \right)^{\top} \right] - \mathbf{f} \right\|_{\Sigma_{\infty}^{m}}^{2},$$

$$= \left\| \mathbb{E} \left[X X^{\top} \right]^{-1} \mathbb{E} \left[X \left(\mathbf{f}^{\top} X + \xi \right)^{\top} \right] - \mathbf{f} \right\|_{\Sigma_{\infty}^{m}}^{2},$$

$$= \left\| \mathbb{E} \left[X X^{\top} \right]^{-1} \mathbb{E} \left[X \left(\mathbf{f}^{\top} X + \xi \right)^{\top} \right] - \mathbf{f} \right\|_{\Sigma_{\infty}^{m}}^{2},$$

$$= \left\| \mathbb{E} \left[X X^{\top} \right]^{-1} \mathbb{E} \left[X \left(\mathbf{f}^{\top} X + \xi \right)^{\top} \right] - \mathbf{f} \right\|_{\Sigma_{\infty}^{m}}^{2},$$

$$= \left\| \mathbb{E} \left[X \mathbb{E} \left[X \right]^{-1} \right] + \mathbb{E} \left[\mathbb{E} \left[X \mathbb{E} \left[X \right]^{-$$

where inequality Eq. (31) holds because \mathbf{D} is non-negative diagonal. Additionally, inequality Eq. (31) holds with equality iff $\mathbf{S}^{-\top}\mathbb{E}\left[X\xi^{\top}\right]$ is in the kernel of \mathbf{D} . Equivalently, iff $\mathbb{E}\left[X\xi^{\top}\right]$ is in kernel of $\mathbf{S}^{\top}\mathbf{D}\mathbf{S} = \mathbf{\Sigma}_{\tilde{Z}}$, which from Lemma 1 holds iff $\mathbb{E}^{\mathfrak{M}}[X|Z] \perp \mathbb{E}^{\mathfrak{M}}[X|\xi]$ a.s.

H.7. Miscellaneous supporting lemmas

Lemma 1 (Gaussian conditional orthogonality lemma) Let $X,Y,Z\in\mathbb{R}^n$ be zero-mean jointly Gaussian random vectors with covariance matrices $\Sigma_X=\mathbb{E}[XX^\top], \ \Sigma_Z=\mathbb{E}[ZZ^\top],$ and cross-covariance $\Sigma_{Y,Z}=\mathbb{E}[YZ^\top].$ Define the conditional expectation

$$\mathbb{E}[Y\mid Z] \coloneqq \left(\mathbb{E}\left[ZZ^{\top}\right]^{-1}\mathbb{E}\left[ZY^{\top}\right]\right)^{\top}Z = \mathbf{\Sigma}_{Y\!,Z}\mathbf{\Sigma}_{Z}^{-1}Z.$$

Then the following are equivalent:

$$X \perp \mathbb{E}[Y \mid Z] = 0$$
 a.s. $\iff \Sigma_X \Sigma_{Y,Z} = 0$.

Proof Since X, Y, Z are jointly Gaussian, $\mathbb{E}[Y \mid Z] = \mathbf{M}Z$ with $\mathbf{M} := \mathbf{\Sigma}_{Y,Z} \mathbf{\Sigma}_Z^{-1}$. The scalar random variable

$$S := X^{\top} \mathbb{E}[Y \mid Z] = X^{\top} \mathbf{M} Z$$

is Gaussian with mean zero. Hence,

$$S = 0$$
 a.s. \iff $Var(S) = 0$.

Compute the variance:

$$\operatorname{Var}(S) = \mathbb{E}\left[S^{2}\right] = \mathbb{E}\left[(X^{\top}\mathbf{M}Z)^{2}\right] = \mathbb{E}\left[Z^{\top}\mathbf{M}^{\top}XX^{\top}\mathbf{M}Z\right].$$

Using independence and zero-mean assumptions,

$$Var(S) = tr(\mathbf{M}^{\top} \mathbf{\Sigma}_X \mathbf{M} \mathbf{\Sigma}_Z).$$

Since covariance matrices are positive semidefinite, Var(S) = 0 iff

$$\mathbf{\Sigma}_X^{1/2}\mathbf{M}\mathbf{\Sigma}_Z^{1/2}=\mathbf{0} \implies \mathbf{\Sigma}_X\mathbf{M}\mathbf{\Sigma}_Z=\mathbf{0}.$$

Substituting $\mathbf{M} = \mathbf{\Sigma}_{Y,Z} \mathbf{\Sigma}_Z^{-1}$ gives

$$\Sigma_X \Sigma_{Y,Z} = \mathbf{0},$$

completing the proof.

Lemma 2 (SPD and PSD simultaneous denationalization via congruence) For any $n \times n$ matrices $\mathbf{A} \succ \mathbf{0}$, $\mathbf{B} \succcurlyeq \mathbf{0}$, there exists an invertible $\mathbf{S} \in \mathbb{R}^{n \times n}$ and non-negative diagonal $\mathbf{D} \in \mathbb{R}^{n \times n}$ such that

$$\mathbf{A} = \mathbf{S}^{\mathsf{T}} \mathbf{S}, \qquad \qquad \mathbf{B} = \mathbf{S}^{\mathsf{T}} \mathbf{D} \mathbf{S}.$$

Proof This is similar to Theorem 7.6.4 in (Horn and Johnson, 1985, p. 465) for two SPD matrices. We proceed similarly; Since **A** is SPD, it admits a unique SPD square root $\mathbf{A}^{1/2}$. Define

$$\mathbf{C} := \mathbf{A}^{-1/2} \mathbf{B} \mathbf{A}^{-1/2}.$$

which is SPD. By the spectral theorem, there exists an orthogonal matrix U such that

$$C = U^{\mathsf{T}}DU$$
,

where **D** is diagonal with non-negative entries (the eigenvalues of **C**). Set

$$\mathbf{S} := \mathbf{U}\mathbf{A}^{1/2}.$$

Then

$$\mathbf{S}^{\mathsf{T}}\mathbf{S} = \mathbf{A}^{1/2}\mathbf{U}^{\mathsf{T}}\mathbf{U}\mathbf{A}^{1/2} = \mathbf{A}^{1/2}\mathbf{I}\mathbf{A}^{1/2} = \mathbf{A},$$

and

$$\mathbf{S}^{\top}\mathbf{D}\mathbf{S} = \mathbf{A}^{1/2}\mathbf{U}^{\top}\mathbf{D}\mathbf{U}\mathbf{A}^{1/2} = \mathbf{A}^{1/2}\mathbf{C}\mathbf{A}^{1/2} = \mathbf{B}.$$

Since $A^{1/2}$ and U are invertible, S is invertible, completing the proof.

Lemma 3 (solvability of simultaneous SEM) The SEM \mathfrak{M} in Example 2 is solvable iff $\mathbf{f}^{\top} \boldsymbol{\tau}^{\top} \neq 1$, in which case the following solution defines the reduced form of the SEM.

$$\begin{bmatrix} X \\ Y \end{bmatrix} = \begin{bmatrix} \mathbf{I}_m & -\boldsymbol{\tau}^\top \\ -\mathbf{f}^\top & 1 \end{bmatrix}^{-1} \left(\begin{bmatrix} \boldsymbol{\Gamma}^\top \\ \mathbf{0}_{1\times k} \end{bmatrix} Z + \begin{bmatrix} \mathbf{T}^\top \\ \boldsymbol{\epsilon}^\top \end{bmatrix} C + \sigma \cdot \begin{bmatrix} N_X \\ N_Y \end{bmatrix} \right),$$

Similarly, SEM \mathfrak{A} in Example 1 solves for $\mathbf{f}^{\top} \boldsymbol{\tau}^{\top} \neq \kappa^{-1}$.

Proof We re-state the SEM $\mathfrak M$ in the following block form

$$\begin{bmatrix} X \\ Y \end{bmatrix} = \begin{bmatrix} \mathbf{0}_{m \times m} & \boldsymbol{\tau}^{\top} \\ \mathbf{f}^{\top} & \mathbf{0}_{1 \times 1} \end{bmatrix} \begin{bmatrix} X \\ Y \end{bmatrix} + \begin{bmatrix} \mathbf{\Gamma}^{\top} \\ \mathbf{0}_{1 \times k} \end{bmatrix} Z + \begin{bmatrix} \mathbf{T}^{\top} \\ \boldsymbol{\epsilon}^{\top} \end{bmatrix} C + \sigma \cdot \begin{bmatrix} N_{X} \\ N_{Y} \end{bmatrix},$$

$$\Rightarrow \begin{bmatrix} \mathbf{I}_{m} & -\boldsymbol{\tau}^{\top} \\ -\mathbf{f}^{\top} & 1 \end{bmatrix} \cdot \begin{bmatrix} X \\ Y \end{bmatrix} = \begin{bmatrix} \mathbf{\Gamma}^{\top} \\ \mathbf{0}_{1 \times k} \end{bmatrix} Z + \begin{bmatrix} \mathbf{T}^{\top} \\ \boldsymbol{\epsilon}^{\top} \end{bmatrix} C + \sigma \cdot \begin{bmatrix} N_{X} \\ N_{Y} \end{bmatrix},$$

solving for (X, Y) involves inverting the block matrix on the LHS. The result immediately follows from Proposition 2.8.7 in (Bernstein, 2009, p. 108), via the Schur complement formula for block matrix inversion.

Lemma 4 (DA and invariance) In SEM $\mathfrak A$ of example from Example 1 iff $\gamma \to \infty$, then

$$\hat{\mathbf{h}}_{DA_G+ERM}^{\mathfrak{A}} \in \operatorname*{argmin}_{\mathbf{h}} R_{DA_G+IV}^{\mathfrak{A}}(\mathbf{h}),$$

Proof We have

$$\hat{\mathbf{h}}_{\mathrm{DA}_{G}+\mathrm{ERM}}^{\mathfrak{A}} = \mathbb{E}\left[(GX)(GX)^{\top} \right]^{-1} \mathbb{E}\left[(GX)Y^{\top} \right], \\
= \mathbb{E}\left[\left(X + \gamma \tilde{G} \right) \left(X + \gamma \tilde{G} \right)^{\top} \right]^{-1} \mathbb{E}\left[\left(X + \gamma \tilde{G} \right) Y^{\top} \right], \quad (\text{Represent } \tilde{G} := \mathbf{\Gamma}^{\top} G.) \\
= \mathbb{E}\left[\left(X + \gamma \tilde{G} \right) \left(X + \gamma \tilde{G} \right)^{\top} \right]^{-1} \mathbb{E}\left[XY^{\top} \right], \quad (\tilde{G} \perp \!\!\!\perp Y \text{ by definition of DA.}) \\
= \mathbb{E}\left[XX^{\top} + \gamma X \tilde{G}^{\top} + \gamma \tilde{G} X^{\top} + \gamma^{2} \tilde{G} \tilde{G}^{\top} \right]^{-1} \mathbb{E}\left[XY^{\top} \right], \\
= \mathbb{E}\left[XX^{\top} + \gamma^{2} \tilde{G} \tilde{G}^{\top} \right]^{-1} \mathbb{E}\left[XY^{\top} \right], \quad (G \text{ independently sampled, } \Rightarrow \tilde{G} \perp \!\!\!\!\perp X.) \\
= \left(\mathbf{\Sigma}_{X} + \gamma^{2} \mathbf{\Sigma}_{\tilde{G}} \right)^{-1} \mathbb{E}\left[XY^{\top} \right], \\
= \left(\mathbf{S}^{\top} \mathbf{S} + \mathbf{S}^{\top} \mathbf{D} \mathbf{S} \right)^{-1} \left(\mathbb{E}\left[XY^{\top} \right] \right), \quad (\text{From Lemma 2.}) \\
= \mathbf{S}^{-1} \left(\mathbf{I}_{m} + \gamma^{2} \mathbf{D} \right)^{-1} \mathbf{S}^{-\top} \left(\mathbb{E}\left[XY^{\top} \right] \right)$$

Now,

$$\nabla_{\mathbf{h}} R_{\mathrm{DA}_{G}+\mathrm{IV}}^{\mathfrak{A}}(\mathbf{h}) = \nabla_{\mathbf{h}} \mathbb{E} \left[\left\| \mathbf{h}^{\top} \mathbb{E}[GX \mid G] - Y \right\|^{2} \right],$$

$$= \mathbb{E} \left[\nabla_{\mathbf{h}} \left\| \mathbf{h}^{\top} \mathbb{E}[GX \mid G] - Y \right\|^{2} \right],$$

$$= \mathbb{E} \left[\mathbb{E}[GX \mid G] \left(\mathbf{h}^{\top} \mathbb{E}[GX \mid G] - Y \right)^{\top} \right],$$

$$= \mathbb{E} \left[\left(\gamma \tilde{G} \right) \left(\gamma \mathbf{h}^{\top} \tilde{G} - Y \right)^{\top} \right], \quad \text{(First stage regression } \mathbb{E}[GX \mid G] = \gamma \tilde{G}.)$$

$$= \gamma^{2} \Sigma_{\tilde{G}} \mathbf{h}. \qquad (\tilde{G} \perp \!\!\!\perp Y \text{ by definition of DA.)}$$

Setting $\nabla_{\mathbf{h}} R^{\mathfrak{A}}_{\mathrm{DA}_G+\mathrm{IV}}(\mathbf{h}) = \gamma^2 \mathbf{\Sigma}_{\tilde{G}} \mathbf{h} = \mathbf{0}_m$, we see that $\Pi^{\mathfrak{A}}_{\mathrm{DA}_G+\mathrm{IV}}$ projects onto the kernel of $\mathbf{\Sigma}_{\tilde{G}}$.

From Eq. (32), we can see that since **D** is a non-negative diagonal, therefore $\hat{\mathbf{h}}_{\mathrm{DA}_G+\mathrm{ERM}}^{\mathfrak{A}}$ only lies in the kernel of $\boldsymbol{\Sigma}_{\tilde{G}} = \mathbf{S}^{\top}\mathbf{D}\mathbf{S}$ when $\gamma \to \infty$. Hence proved that only when $\gamma \to \infty$,

$$\hat{\mathbf{h}}_{\mathrm{DA}_G + \mathrm{ERM}}^{\mathfrak{A}} \in \operatorname*{argmin}_{\mathbf{h}} R_{\mathrm{DA}_G + \mathrm{IV}}^{\mathfrak{A}}(\mathbf{h}).$$

42