

Counterfactual Cocycles: Noise Invariant and Coherent Transport-based Couplings

Hugh Dance¹, Benjamin Bloem-Reddy²

Causal Modelling and Inference Annual Workshop,
CAUSALI-T-AI, Institut Henri Poincaré

August 25, 2025

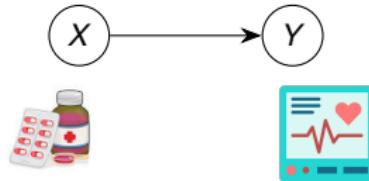


¹PhD Student, Gatsby Unit (P. Orbanz lab), UCL

²Assistant Professor, Department of Statistics, University of British Columbia

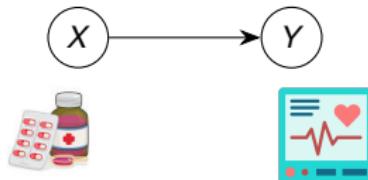
RCT as the “gold-standard” for marginal counterfactual inference

Goal: Assess effectiveness of medication $X \in \{0, 1\}$ on symptoms $Y := (Y_1, \dots, Y_p)$.



RCT as the “gold-standard” for marginal counterfactual inference

Goal: Assess effectiveness of medication $X \in \{0, 1\}$ on symptoms $Y := (Y_1, \dots, Y_p)$.

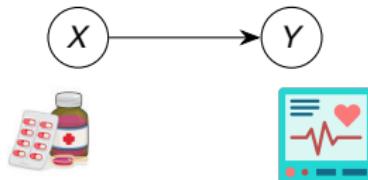


Randomized Control Trial (RCT)

- Observe treatment units $\{Y^{(i)}(1)\}_{i=1}^{n_1}$ and control units $\{Y^{(i)}(0)\}_{i=n_1+1}^{n_1+n_0}$.
- Potential outcomes satisfy $Y^{(i)}(x) \sim \mathbb{P}_{Y|X=x}$

RCT as the “gold-standard” for marginal counterfactual inference

Goal: Assess effectiveness of medication $X \in \{0, 1\}$ on symptoms $Y := (Y_1, \dots, Y_p)$.



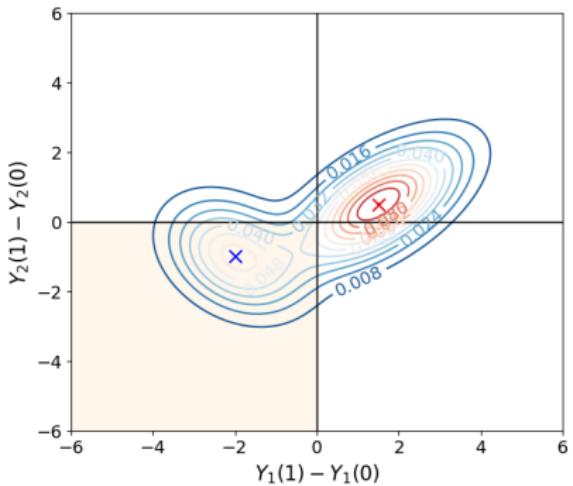
Randomized Control Trial (RCT)

- Observe treatment units $\{Y^{(i)}(1)\}_{i=1}^{n_1}$ and control units $\{Y^{(i)}(0)\}_{i=n_1+1}^{n_1+n_0}$.
- Potential outcomes satisfy $Y^{(i)}(x) \sim \mathbb{P}_{Y|X=x}$

$$\text{ATE} = \mathbb{E}[Y(1)] - \mathbb{E}[Y(0)] \implies \widehat{\text{ATE}} = \frac{1}{n_1} \sum_{i=1}^{n_1} Y^{(i)}(1) - \frac{1}{n_0} \sum_{i=n_1+1}^{n_0} Y^{(i)}(0)$$

When is the “gold-standard” not enough?

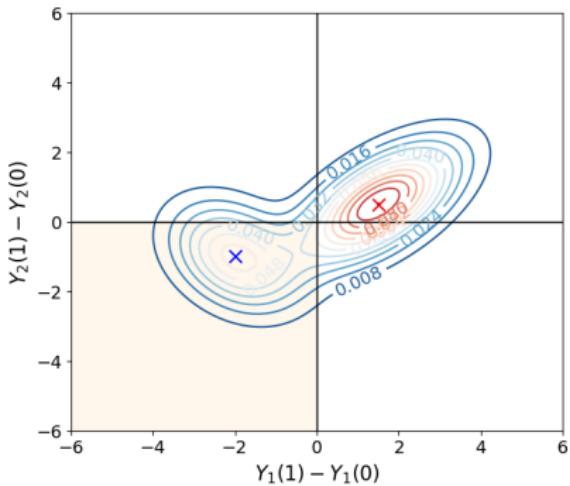
Goal: Quantify possible treatment harms



³Shen et al. (2013). Treatment Benefit and Treatment Harm Rate to Characterize Heterogeneity in Treatment Effect, *Biometrics*.

When is the “gold-standard” not enough?

Goal: Quantify possible treatment harms

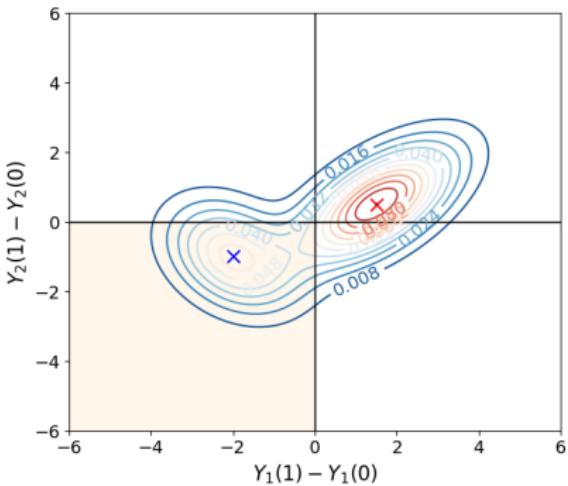


Treatment Harm Rate (THR)³ := $\mathbb{P}(Y(1) \prec Y(0)) = \mathbb{P}(\{Y_1(1) \prec Y_1(0)\} \cap \{Y_2(1) - Y_2(0)\})$

³Shen et al. (2013). Treatment Benefit and Treatment Harm Rate to Characterize Heterogeneity in Treatment Effect, *Biometrics*.

When is the “gold-standard” not enough?

Goal: Quantify possible treatment harms



Treatment Harm Rate (THR)³ := $\mathbb{P}(Y(1) \prec Y(0)) = \mathbb{P}(\{Y_1(1) \prec Y_1(0)\} \cap \{Y_2(1) - Y_2(0)\})$

Issue: Only observe $Y^{(i)}(0)$ or $Y^{(i)}(1)$ for each unit - cannot estimate THR from data!

³Shen et al. (2013). Treatment Benefit and Treatment Harm Rate to Characterize Heterogeneity in Treatment Effect, *Biometrics*.

The Need for Counterfactual Couplings

Underlying issue: need coupling $\mathbb{P}_{Y(1), Y(0)}$ to estimate $\mathbb{P}_{Y(1) - Y(0)}$

$$\text{THR} := \iint \mathbf{1}\{y_1 - y_0 \leq 0\} d\mathbb{P}_{Y(1), Y(0)}(y_1, y_0)$$

⁴Lee, M.J. (2000). Median Treatment Effect in Randomized Trials, JRSSB-B.

⁵Kallus, N. (2023). Treatment effect risk: Bounds and inference, Management Science.

⁶Shen et al. (2013). Treatment Benefit and Treatment Harm Rate to Characterize Heterogeneity in Treatment Effect, Biometrics.

The Need for Counterfactual Couplings

Underlying issue: need coupling $\mathbb{P}_{Y(1), Y(0)}$ to estimate $\mathbb{P}_{Y(1) - Y(0)}$

$$\text{THR} := \iint \mathbf{1}\{y_1 - y_0 \leq 0\} d\mathbb{P}_{Y(1), Y(0)}(y_1, y_0)$$

Other Examples

- Median Treatment Effect (MTE)⁴: $\text{Med}(Y(1) - Y(0))$
- Conditional Value at Risk (CVar)⁵: $\mathbb{E}[Y(1) - Y(0)|Y(1) - Y(0) \leq q_\alpha]$
- Treatment Benefit Rate (TBR)⁶: $\mathbb{P}(Y(1) \succ Y(0))$

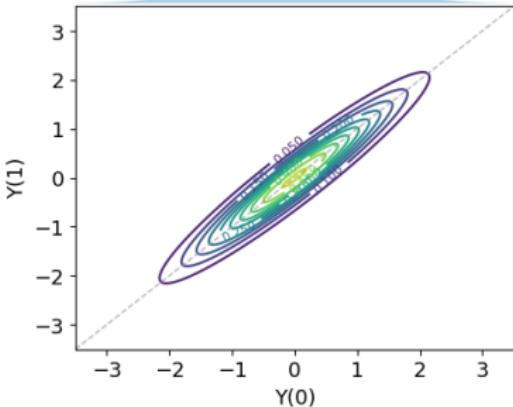
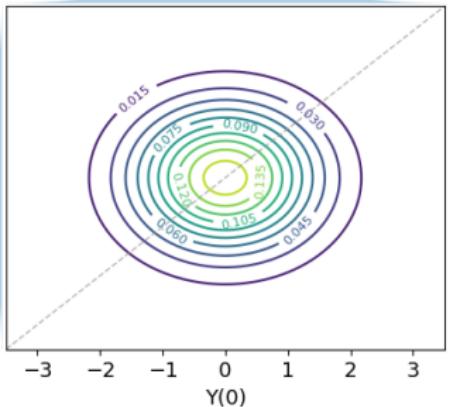
⁴Lee, M.J. (2000). Median Treatment Effect in Randomized Trials, JRSSB-B.

⁵Kallus, N. (2023). Treatment effect risk: Bounds and inference, Management Science.

⁶Shen et al. (2013). Treatment Benefit and Treatment Harm Rate to Characterize Heterogeneity in Treatment Effect, Biometrics.

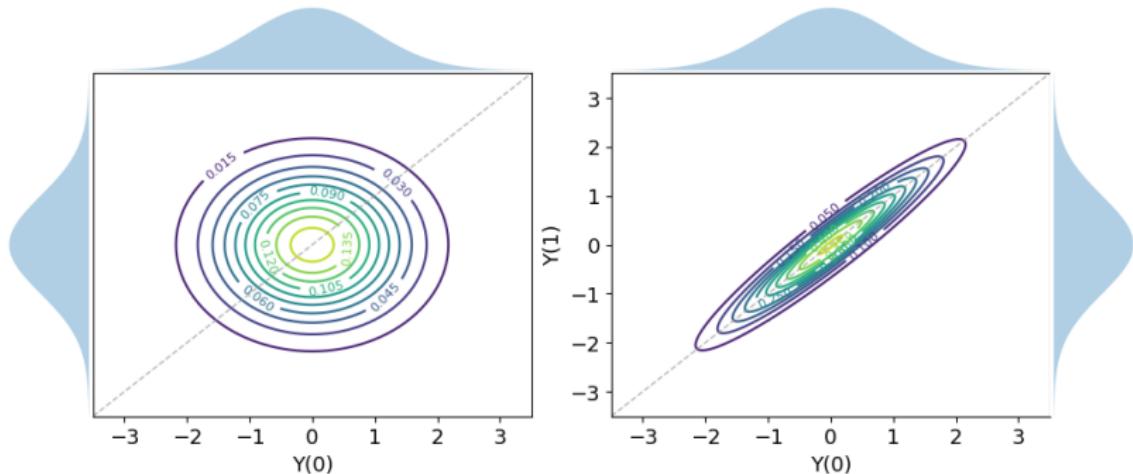
Counterfactual Couplings: Identification Challenge

- Infinitely many admissible couplings $\pi \in \Pi(\mathbb{P}_{Y(1)}, \mathbb{P}_{Y(0)})$
- Can never observe joint samples $(Y(1), Y(0))$



Counterfactual Couplings: Identification Challenge

- Infinitely many admissible couplings $\pi \in \Pi(\mathbb{P}_{Y(1)}, \mathbb{P}_{Y(0)})$
- Can never observe joint samples $(Y(1), Y(0))$

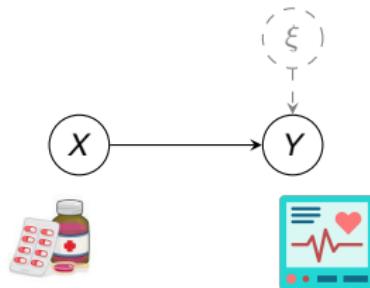


Existing Approaches

- Structural Causal Models (SCMs)
- Optimal Transport Methods (OT)

Limitations of Existing Methods

Structural Causal Models⁷

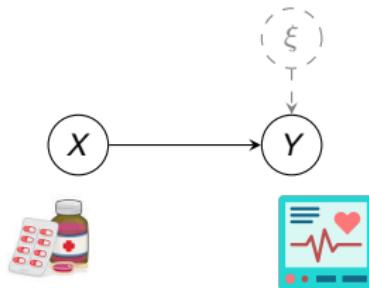


Structural Causal Model (SCM): $Y = f(X, \xi)$, $\xi \sim \mathbb{P}_\xi$, $\xi \perp\!\!\!\perp X$

Counterfactuals: $Y(x) = f(x, \xi) \implies$ induced coupling: $(Y(1), Y(0)) = (f_1(\xi), f_0(\xi))$

⁷Spirites et al. (2000), Pearl et al. (2009), Bongers et al. (2021).

Structural Causal Models⁷



Structural Causal Model (SCM): $Y = f(X, \xi)$, $\xi \sim \mathbb{P}_\xi$, $\xi \perp\!\!\!\perp X$

Counterfactuals: $Y(x) = f(x, \xi) \implies$ induced coupling: $(Y(1), Y(0)) = (f_1(\xi), f_0(\xi))$

... how to identify (f, \mathbb{P}_ξ) from data?

⁷Spirites et al. (2000), Pearl et al. (2009), Bongers et al. (2021).

Structural Causal Model Non-Identifiability

Transport Problem : find (f, \mathbb{P}_ξ) such that $(f_x)_\# : \mathbb{P}_\xi \mapsto \mathbb{P}_{Y|X=x}$ (i.e. $f_x(\xi) \sim \mathbb{P}_{Y|X=x}$)

Structural Causal Model Non-Identifiability

Transport Problem : find (f, \mathbb{P}_ξ) such that $(f_x)_\# : \mathbb{P}_\xi \mapsto \mathbb{P}_{Y|X=x}$ (i.e. $f_x(\xi) \sim \mathbb{P}_{Y|X=x}$)

Issue: Can find $f^{(1)}, f^{(2)}$ such that $f_x^{(1)}(\xi) =_d f_x^{(2)}(\xi) \sim \mathbb{P}_{Y|X=x}$ for each x but:

$$(f_1^{(1)}(\xi), f_0^{(1)}(\xi)) \neq_d (f_1^{(2)}(\xi), f_0^{(2)}(\xi))$$

Structural Causal Model Non-Identifiability

Transport Problem : find (f, \mathbb{P}_ξ) such that $(f_x)_\# : \mathbb{P}_\xi \mapsto \mathbb{P}_{Y|X=x}$ (i.e. $f_x(\xi) \sim \mathbb{P}_{Y|X=x}$)

Issue: Can find $f^{(1)}, f^{(2)}$ such that $f_x^{(1)}(\xi) =_d f_x^{(2)}(\xi) \sim \mathbb{P}_{Y|X=x}$ for each x but:

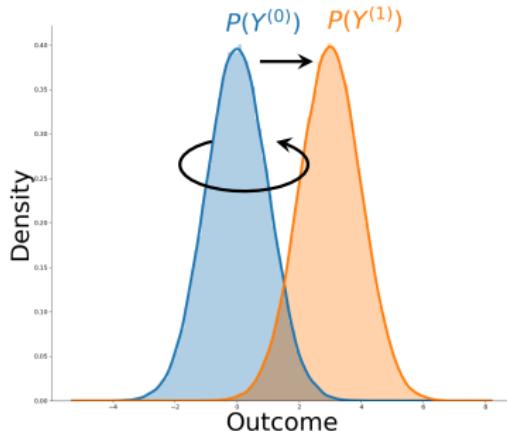
$$(f_1^{(1)}(\xi), f_0^{(1)}(\xi)) \neq_d (f_1^{(2)}(\xi), f_0^{(2)}(\xi))$$

Example:

$$\xi \sim \mathcal{N}(0, 1)$$

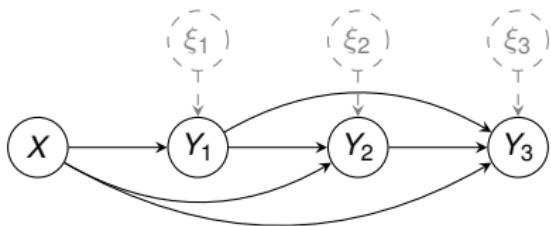
$$f^{(1)}(x, \xi) = \beta x + \xi$$

$$f^{(2)}(x, \xi) = \beta x + (2x - 1)\xi$$



Couplings don't agree! $(\xi, 1 + \xi) \neq_d (\xi, 1 - \xi)$

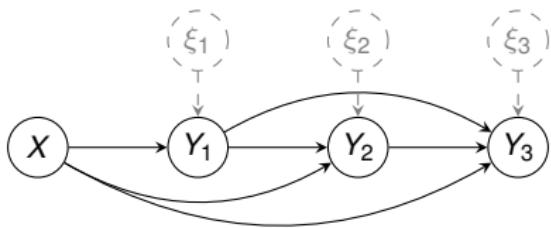
Identifiability via Bijective Causal Models (BCMs)



Causal Ordering: $Y_1 \prec Y_2 \prec \dots \prec Y_p \quad \Rightarrow \quad \text{Acyclic-SCM: } Y_j = f_j(X, Y_{<j}, \xi_j) \quad \forall j$

⁸Bogachev, V.I. et al. (2005) "Triangular transformations of measures", Sbornik: Mathematics.

Identifiability via Bijective Causal Models (BCMs)

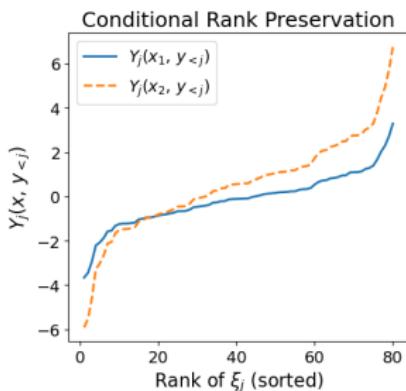


Causal Ordering: $Y_1 \prec Y_2 \prec \dots \prec Y_p \implies$ **Acyclic-SCM:** $Y_j = f_j(X, Y_{<j}, \xi_j) \quad \forall j$

Conditions for Identifiability⁸

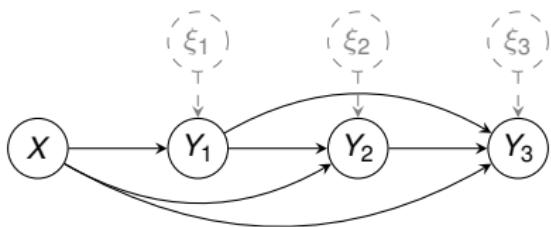
- Each f_j is bijective on ξ_j
- $\mathbb{P}_{Y|X=x}$ and \mathbb{P}_ξ are abs. continuous on \mathbb{R}^d

$\Rightarrow f_X : \mathcal{E} \rightarrow \mathbb{Y}$ is *triangular, monotone, increasing*



⁸Bogachev, V.I. et al. (2005) "Triangular transformations of measures", Sbornik: Mathematics.

Identifiability via Bijective Causal Models (BCMs)

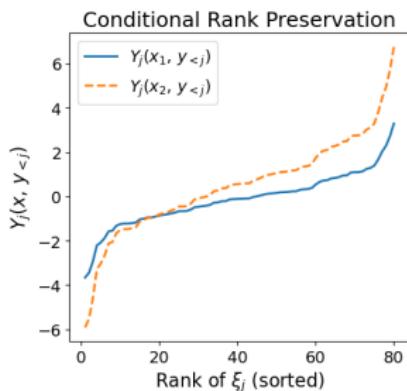


Causal Ordering: $Y_1 \prec Y_2 \prec \dots \prec Y_p \implies$ **Acyclic-SCM:** $Y_j = f_j(X, Y_{<j}, \xi_j) \quad \forall j$

Conditions for Identifiability⁸

- Each f_j is bijective on ξ_j
- $\mathbb{P}_{Y|X=x}$ and \mathbb{P}_ξ are abs. continuous on \mathbb{R}^d

$\Rightarrow f_X : \mathcal{E} \rightarrow \mathbb{Y}$ is *triangular, monotone, increasing*



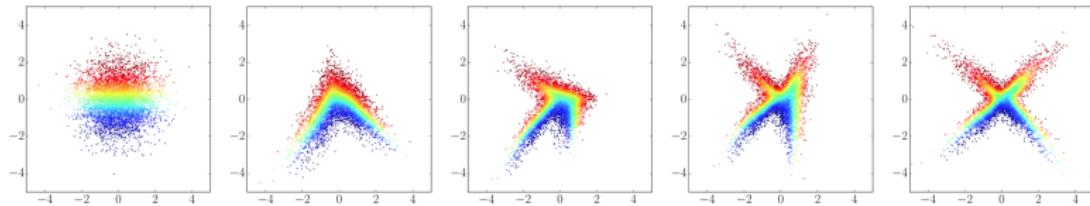
Can augment model with measured causes Z (age, gender, gene expression...)

⁸Bogachev, V.I. et al. (2005) "Triangular transformations of measures", Sbornik: Mathematics.

How to model BCMs?

Popular Approach⁹: Fix ‘base’ distribution \mathbb{P}_ξ and learn diffeomorphisms $(f_x)_{x \in \mathbb{X}}$

$$y(x) = f_x^{(L)} \circ f_x^{(L-1)} \circ \cdots \circ f_x^{(1)}(\xi) \implies \log p(y|x) = \log p_\xi(f_x^{-1}(y)) + \log |\nabla_y f_x^{-1}(y)|$$

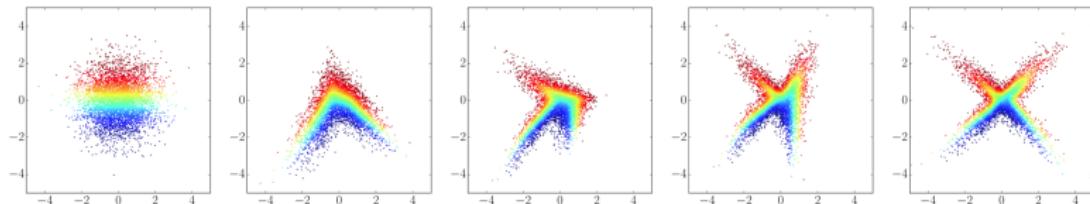


⁹Khemekhem et al. AISTATS'21, Nasr-Esfahany et al. ICML'22, Javaloy et al. NeurIPS'24, Zhou et al. AAAI'25

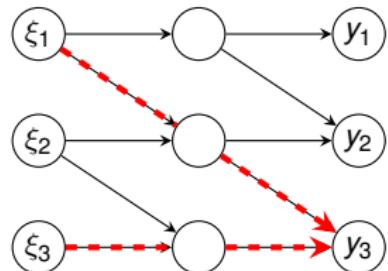
How to model BCMs?

Popular Approach⁹: Fix ‘base’ distribution \mathbb{P}_ξ and learn diffeomorphisms $(f_x)_{x \in \mathbb{X}}$

$$y(x) = f_x^{(L)} \circ f_x^{(L-1)} \circ \cdots \circ f_x^{(1)}(\xi) \implies \log p(y|x) = \log p_\xi(f_x^{-1}(y)) + \log |\nabla_y f_x^{-1}(y)|$$



Autoregressive Normalizing Flows: Respect causal ordering when stacked!

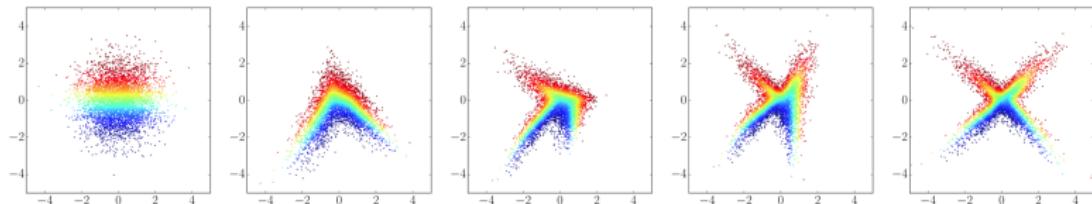


⁹Khemekhem et al. AISTATS'21, Nasr-Esfahany et al. ICML'22, Javaloy et al. NeurIPS'24, Zhou et al. AAAI'25

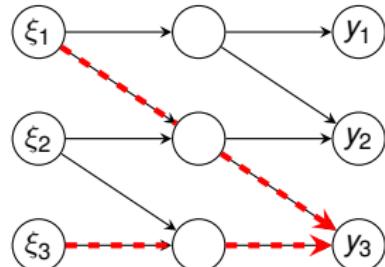
How to model BCMs?

Popular Approach⁹: Fix ‘base’ distribution \mathbb{P}_ξ and learn diffeomorphisms $(f_x)_{x \in \mathbb{X}}$

$$y(x) = f_x^{(L)} \circ f_x^{(L-1)} \circ \dots \circ f_x^{(1)}(\xi) \implies \log p(y|x) = \log p_\xi(f_x^{-1}(y)) + \log |\nabla_y f_x^{-1}(y)|$$



Autoregressive Normalizing Flows: Respect causal ordering when stacked!



Model **Autoregressive transform** $f_j(v_{<j}, \xi_j)$

NICE (Additive) $\xi_j + \mu_j(v_{<j})$

MAF (Affine) $\xi_j \mapsto \exp(\lambda_j(v_{<j})) \xi_j + \mu_j(v_{<j})$

NAF (INN) $\xi_j \mapsto \sigma^{-1}(w(v_{<j}) \cdot \sigma(\sigma_j(v_{<j}) \xi_j + \mu_j(v_{<j})))$

NSF (Spline) $\xi_j \mapsto v_j 1_{v_j \notin [-B, B]} + M_j(\xi_j; v_{<j}) 1_{v_j \in [-B, B]}$

⁹Khemekhem et al. AISTATS'21, Nasr-Esfahany et al. ICML'22, Javaloy et al. NeurIPS'24, Zhou et al. AAAI'25

Tail Mis-Specification Problems

- If tail decay of \mathbb{P}_ξ doesn't match $\mathbb{P}_{Y|X=x}$, no bi-lipschitz transport f_x exists!¹⁰
- Heavy tailed $\mathbb{P}_{Y|X=x}$, Gaussian $\hat{\mathbb{P}}_\xi$ = undefined likelihood:

$$|\mathbb{E} \log \hat{p}(Y(x))| \succeq \mathbb{E} \|f_x^{-1}(Y(x))\|^2 = \infty$$

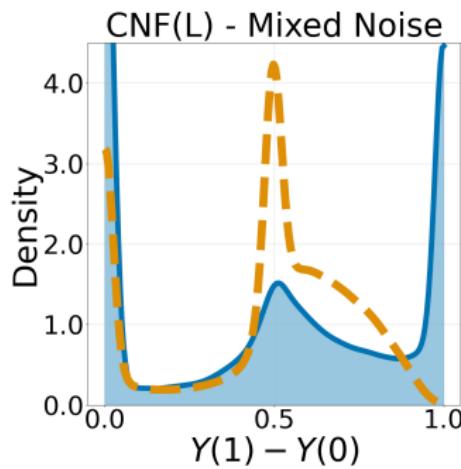
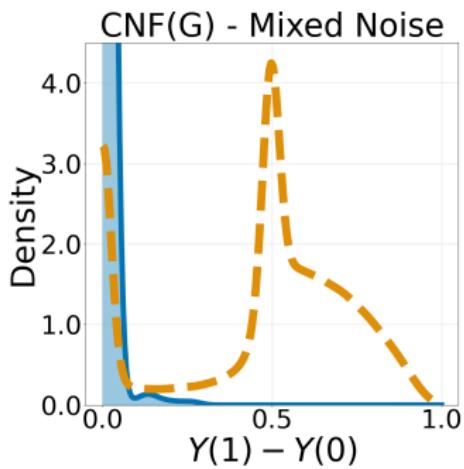
¹⁰Jaini et al. ICML'20, Liang et al. ICML'22

Tail Mis-Specification Problems

- If tail decay of \mathbb{P}_ξ doesn't match $\mathbb{P}_{Y|X=x}$, no bi-lipschitz transport f_x exists!¹⁰
- Heavy tailed $\mathbb{P}_{Y|X=x}$, Gaussian $\hat{\mathbb{P}}_\xi$ = undefined likelihood:

$$|\mathbb{E} \log \hat{p}(Y(x))| \succeq \mathbb{E} \|f_x^{-1}(Y(x))\|^2 = \infty$$

Example : $Y = (X + 1)\xi$, $X \sim \text{Bern}(1/2)$, $\xi \sim \frac{1}{2}|\mathcal{N}(0, 1)| - \frac{1}{2}|NBP(0.1, 0.1)|$



¹⁰Jaini et al. ICML'20, Liang et al. ICML'22

Support Mis-Specification Problems

- Need diffeomorphic $\text{supp}(\hat{\mathbb{P}}_\xi)$ and $\text{supp}(\mathbb{P}_{Y|X=x})$ for existence of $f_x \in \mathcal{F}$
- Likelihood non-identifiability if $\text{supp}((f_x^{-1})_*(\mathbb{P}_{Y|X=x})) \not\supseteq \text{supp}(\hat{\mathbb{P}}_\xi)$:

Support Mis-Specification Problems

- Need diffeomorphic $\text{supp}(\hat{\mathbb{P}}_\xi)$ and $\text{supp}(\mathbb{P}_{Y|X=x})$ for existence of $f_x \in \mathcal{F}$
- Likelihood non-identifiability if $\text{supp}((f_x^{-1})_{\#}(\mathbb{P}_{Y|X=x})) \not\supset \text{supp}(\hat{\mathbb{P}}_\xi)$:

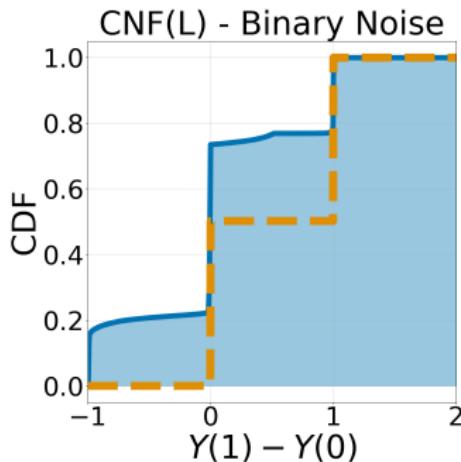
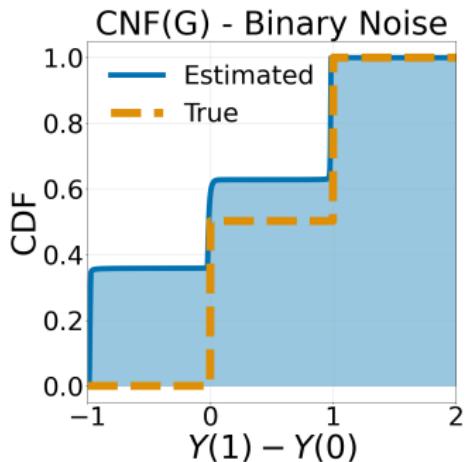
$$\mathbb{E} \log \hat{p}_f(Y(x)) = \int (\log \hat{p}_\xi(f_x^{-1}(y)) + \log |\nabla_y f_x^{-1}(y)|) \mathbb{P}_{Y|X=x}(dy)$$

Support Mis-Specification Problems

- Need diffeomorphic $\text{supp}(\hat{\mathbb{P}}_\xi)$ and $\text{supp}(\mathbb{P}_{Y|X=x})$ for existence of $f_x \in \mathcal{F}$
- Likelihood non-identifiability if $\text{supp}((f_x^{-1})_{\#}(\mathbb{P}_{Y|X=x})) \not\supset \text{supp}(\hat{\mathbb{P}}_\xi)$:

$$\mathbb{E} \log \hat{p}_f(Y(x)) = \int (\log \hat{p}_\xi(f_x^{-1}(y)) + \log |\nabla_y f_x^{-1}(y)|) \mathbb{P}_{Y|X=x}(dy)$$

Example : $Y = (X + 1)\xi$, $X \sim \text{Bern}(1/2)$, $\xi \sim \text{Rad}(1/2)$



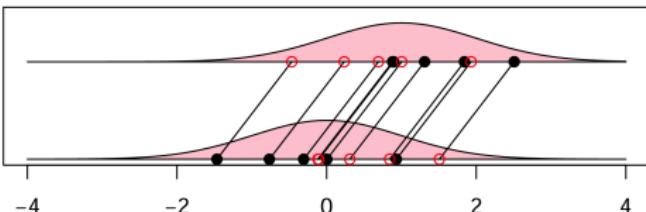
Optimal Transport Methods: An Alternative

Idea: Choose coupling between $Y(1)$, $Y(0)$ using *optimal-transport*:

$$\pi^* = \arg \min_{\pi \in \Pi(\mathbb{P}_{Y(1)}, \mathbb{P}_{Y(0)})} \iint c(y(1), y(0)) d\pi(y(1), y(0))$$

Motivations:

- Conservativism¹
- Counterfactual Similarity²
- Optimal Matching³



¹Balakrishnan et al. (2025) "Conservative inference for counterfactuals." Journal of Causal Inference.

²De Lara et al. (2024) "Transport-based Counterfactual Models" JMLR.

³Charpentier et al (2023) "Optimal transport for counterfactual estimation: A method for causal inference.", Optimal transport statistics for economics and related topics.

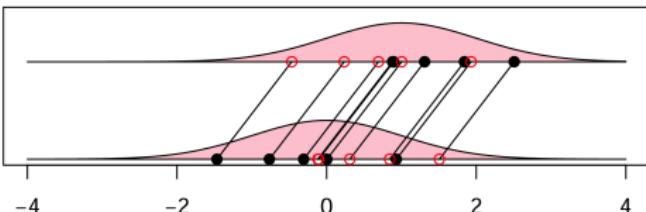
Optimal Transport Methods: An Alternative

Idea: Choose coupling between $Y(1)$, $Y(0)$ using *optimal-transport*:

$$\pi^* = \arg \min_{\pi \in \Pi(\mathbb{P}_{Y(1)}, \mathbb{P}_{Y(0)})} \iint c(y(1), y(0)) d\pi(y(1), y(0))$$

Motivations:

- Conservativism¹
- Counterfactual Similarity²
- Optimal Matching³



Deterministic coupling restriction: $Y(1) = T_{1,0}(Y(0))$

¹Balakrishnan et al. (2025) "Conservative inference for counterfactuals." Journal of Causal Inference.

²De Lara et al. (2024) "Transport-based Counterfactual Models" JMLR.

³Charpentier et al (2023) "Optimal transport for counterfactual estimation: A method for causal inference.", Optimal transport statistics for economics and related topics.

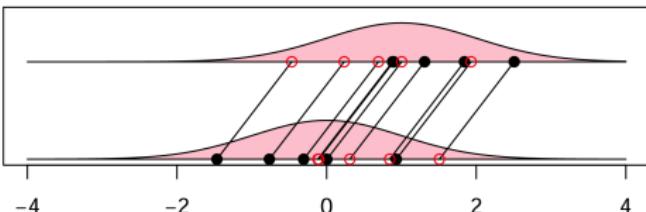
Optimal Transport Methods: An Alternative

Idea: Choose coupling between $Y(1)$, $Y(0)$ using *optimal-transport*:

$$\pi^* = \arg \min_{\pi \in \Pi(\mathbb{P}_{Y(1)}, \mathbb{P}_{Y(0)})} \iint c(y(1), y(0)) d\pi(y(1), y(0))$$

Motivations:

- Conservativism¹
- Counterfactual Similarity²
- Optimal Matching³



Deterministic coupling restriction: $Y(1) = T_{1,0}(Y(0))$

Brenier Maps $T_{1,0}^* = \arg \min_{T \# \mathbb{P}_{Y(0)} = \mathbb{P}_{Y(1)}} \int \|y(0) - T_{1,0}(y(0))\|^2 d\mathbb{P}_{Y(0)}$

¹Balakrishnan et al. (2025) "Conservative inference for counterfactuals." Journal of Causal Inference.

²De Lara et al. (2024) "Transport-based Counterfactual Models" JMLR.

³Charpentier et al (2023) "Optimal transport for counterfactual estimation: A method for causal inference.", Optimal transport statistics for economics and related topics.

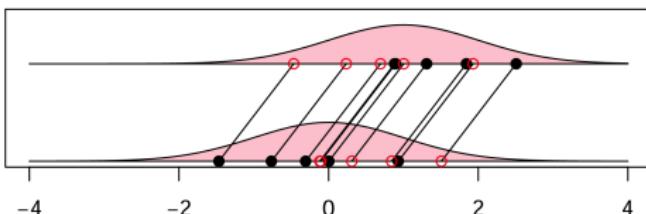
Optimal Transport Methods: An Alternative

Idea: Choose coupling between $Y(1)$, $Y(0)$ using *optimal-transport*:

$$\pi^* = \arg \min_{\pi \in \Pi(\mathbb{P}_{Y(1)}, \mathbb{P}_{Y(0)})} \iint c(y(1), y(0)) d\pi(y(1), y(0))$$

Motivations:

- Conservativism¹
- Counterfactual Similarity²
- Optimal Matching³



Deterministic coupling restriction: $Y(1) = T_{1,0}(Y(0))$

Brenier Maps $T_{1,0}^* = \arg \min_{T \# \mathbb{P}_{Y(0)} = \mathbb{P}_{Y(1)}} \int \|y(0) - T_{1,0}(y(0))\|^2 d\mathbb{P}_{Y(0)}$

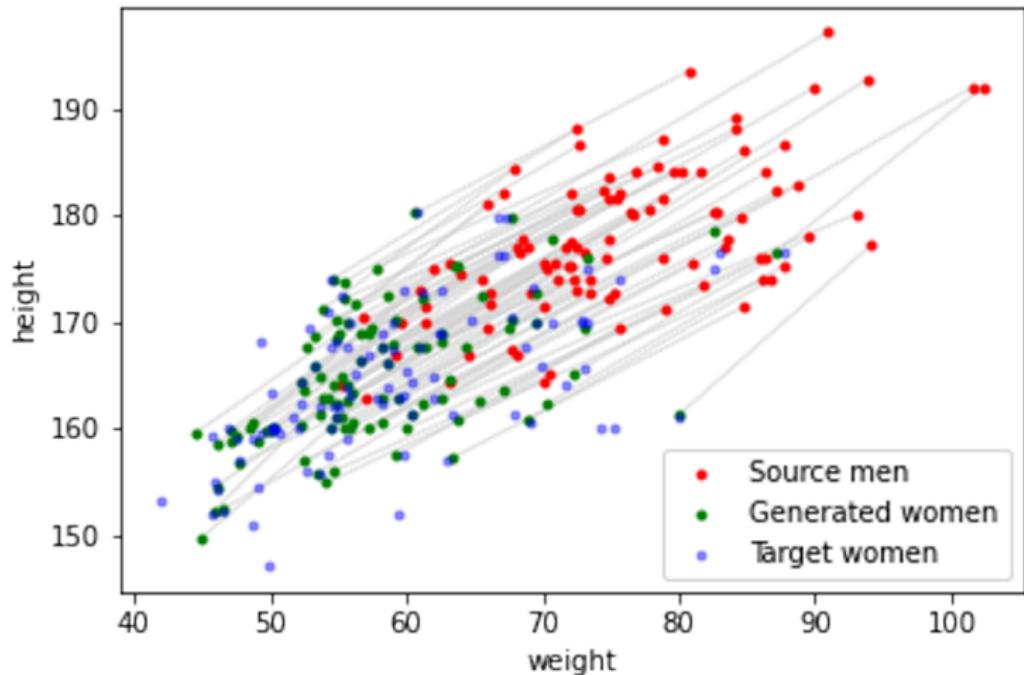
Guarantees identifiability of $T_{1,0}$ between abs. continuous distributions!

¹Balakrishnan et al. (2025) "Conservative inference for counterfactuals." Journal of Causal Inference.

²De Lara et al. (2024) "Transport-based Counterfactual Models" JMLR.

³Charpentier et al (2023) "Optimal transport for counterfactual estimation: A method for causal inference.", Optimal transport statistics for economics and related topics.

Example:¹¹



¹⁰De Lara, L., Gonzales-Sans,A., Asher,N., Risser, L., Loubes, J.M. (2024) "Transport-based Counterfactual Models" JMLR.

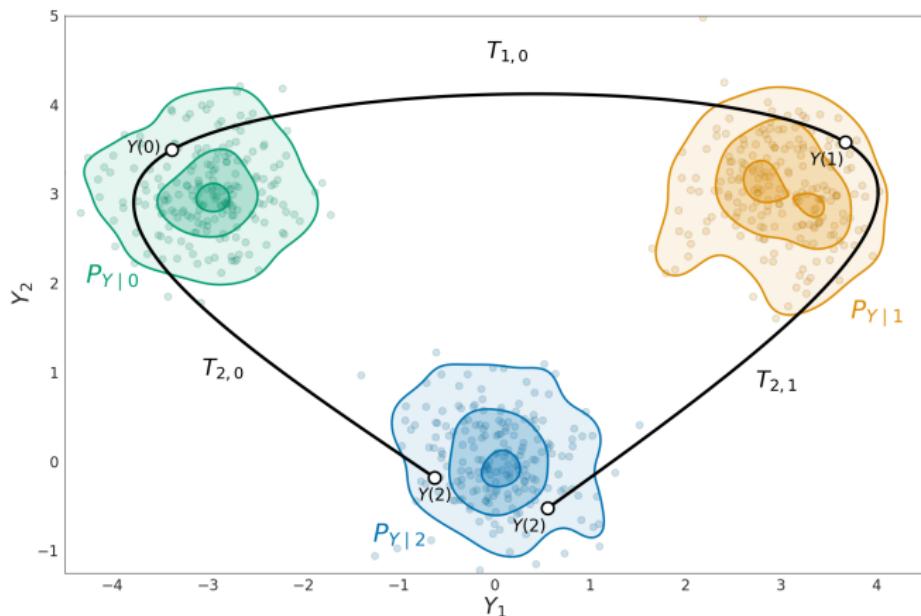
Incoherence of Transports Under Multiple treatments

Issue: If $x \in \{0, 1, 2\}$ and Y multivariate, OT maps not closed under composition!

$$T_{2,1} \circ T_{1,0} \neq T_{2,0}$$

Logical Impossibility:

$$Y(2) = T_{2,1} \circ T_{1,0}(Y(0)) \neq T_{2,0}(Y(0)) = T_{2,0} \circ T_{0,2}(Y(2)) = Y(2)$$



Example: Gaussian Transport

Three Multivariate Gaussians:

$$\mathbb{P}_0 = \mathcal{N}(\mu_0, \Sigma_0), \quad \mathbb{P}_1 = \mathcal{N}(\mu_1, \Sigma_1), \quad \mathbb{P}_2 = \mathcal{N}(\mu_2, \Sigma_2)$$

Brenier (OT) Map:

$$T_{x,x'}(y) = \mu_{x'} + \Sigma_x^{-1/2} (\Sigma_x^{1/2} \Sigma_{x'} \Sigma_x^{1/2})^{1/2} \Sigma_x^{-1/2} (y - \mu_x)$$

Example: Gaussian Transport

Three Multivariate Gaussians:

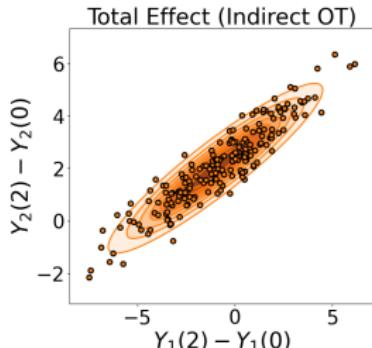
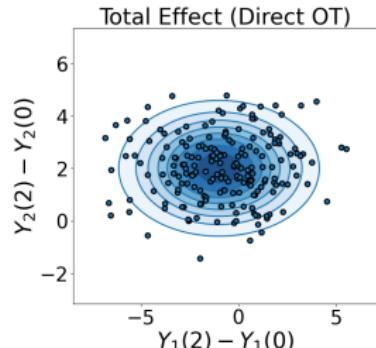
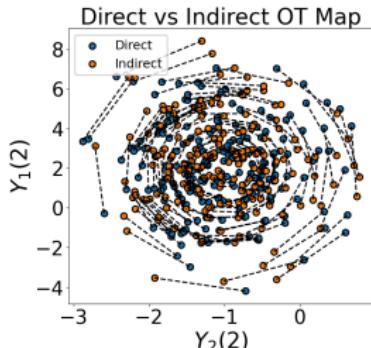
$$\mathbb{P}_0 = \mathcal{N}(\mu_0, \Sigma_0), \quad \mathbb{P}_1 = \mathcal{N}(\mu_1, \Sigma_1), \quad \mathbb{P}_2 = \mathcal{N}(\mu_2, \Sigma_2)$$

Brenier (OT) Map:

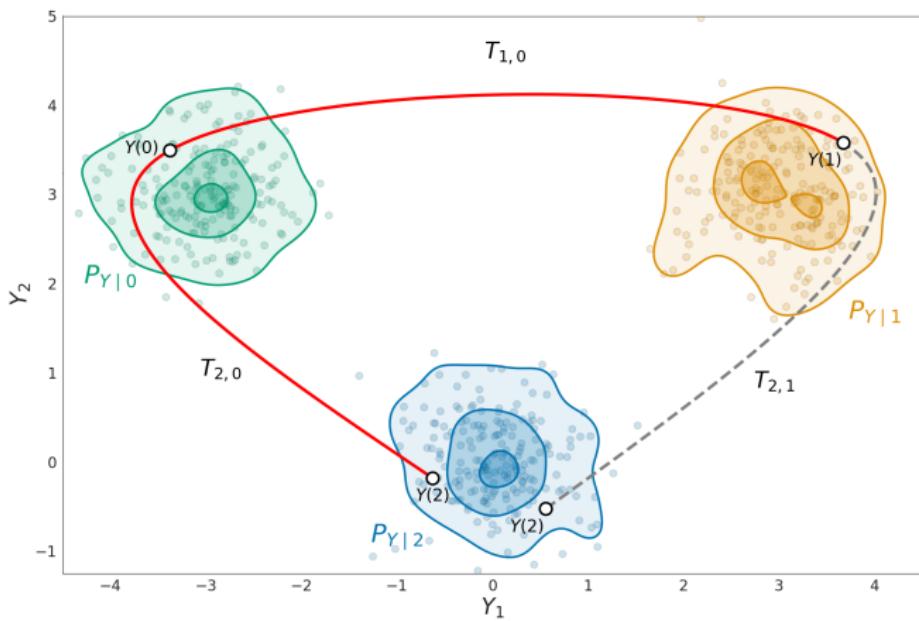
$$T_{x,x'}(y) = \mu_{x'} + \Sigma_x^{-1/2} (\Sigma_x^{1/2} \Sigma_{x'} \Sigma_x^{1/2})^{1/2} \Sigma_x^{-1/2} (y - \mu_x)$$

Illustration: Draw $Y(0) \sim \mathbb{P}_0$ and impute $Y(1)$, $Y(2)$

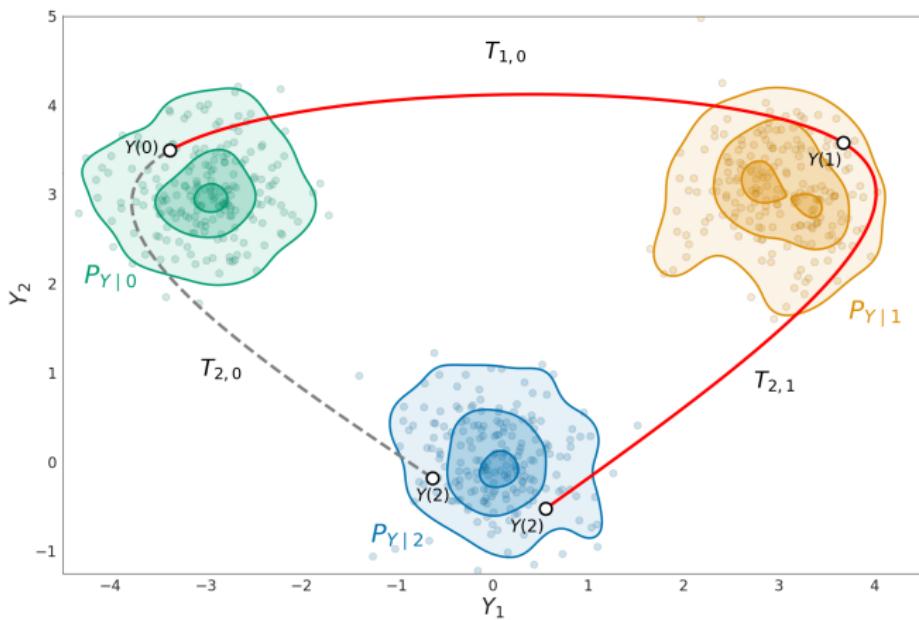
- **Direct:** $Y(1) = T_{1,0}(Y(0))$, $Y(2) = T_{2,0}(Y(0))$
- **Indirect:** $Y(1) = T_{1,0}(Y(0))$, $Y(2) = T_{2,1} \circ T_{1,0}(Y(0))$



Selecting Transport Subsets Induces Coupling Non-Identifiability



Selecting Transport Subsets Induces Coupling Non-Identifiability



Transport-based Models with Cocycles

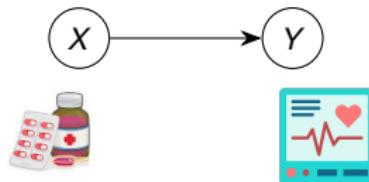
Goals

Modeling and Estimation Framework for counterfactual couplings that:

- Guarantees Coherence
- Avoids Latent Noise Assumptions

Formal Set-up

Treatment $X \in \mathbb{R}$, Outcomes $Y := (Y_1, \dots, Y_p) \in \mathbb{R}^p$

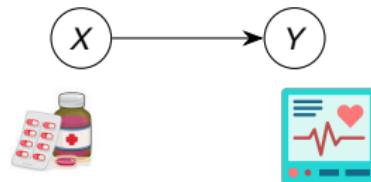


Counterfactuals: There exist 'potential outcomes' $\{Y(x)\}_{x \in \mathbb{X}}$ satisfying

- Consistency: $X = x \implies Y(x) = Y$
- Exchangeability: $Y(x) \perp\!\!\!\perp X$

Formal Set-up

Treatment $X \in \mathbb{R}$, Outcomes $Y := (Y_1, \dots, Y_p) \in \mathbb{R}^p$



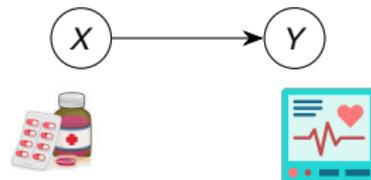
Counterfactuals: There exist ‘potential outcomes’ $\{Y(x)\}_{x \in \mathbb{X}}$ satisfying

- Consistency: $X = x \implies Y(x) = Y$
- Exchangeability: $Y(x) \perp\!\!\!\perp X$

$$\implies Y(x) \sim \mathbb{P}_{Y|X=x}$$

Formal Set-up

Treatment $X \in \mathbb{R}$, Outcomes $Y := (Y_1, \dots, Y_p) \in \mathbb{R}^p$



Counterfactuals: There exist ‘potential outcomes’ $\{Y(x)\}_{x \in \mathbb{X}}$ satisfying

- Consistency: $X = x \implies Y(x) = Y$
- Exchangeability: $Y(x) \perp\!\!\!\perp X$

$$\implies Y(x) \sim \mathbb{P}_{Y|X=x}$$

Transport-based model:

$$Y(x) = T_{x,x'}(Y(x')), \quad \forall x, x' \in \mathbb{X} \quad (\text{CC})$$

Necessary Algebraic Properties

$$Y(x) = T_{x,x'}(Y(x')), \quad \forall x, x' \in \mathbb{X} \quad (\text{CC})$$

Necessary Algebraic Properties

$$Y(x) = T_{x,x'}(Y(x')), \quad \forall x, x' \in \mathbb{X} \quad (\text{CC})$$

1. **Distribution Adaptedness:** For each $x, x' \in \mathbb{X}$,

$$(T_{x',x})_\# \mathbb{P}_{Y|X=x} = \mathbb{P}_{Y|X=x'} \quad (\text{DA})$$

2. **Identity:** For each $x \in \mathbb{X}$:

$$T_{x,x}(y) = y, \quad \forall y \in \mathbb{Y}_x \quad (\text{ID})$$

Necessary Algebraic Properties

$$Y(x) = T_{x,x'}(Y(x')), \quad \forall x, x' \in \mathbb{X} \quad (\text{CC})$$

1. **Distribution Adaptedness:** For each $x, x' \in \mathbb{X}$,

$$(T_{x',x})_{\#} \mathbb{P}_{Y|X=x} = \mathbb{P}_{Y|X=x'} \quad (\text{DA})$$

2. **Identity:** For each $x \in \mathbb{X}$:

$$T_{x,x}(y) = y, \quad \forall y \in \mathbb{Y}_x \quad (\text{ID})$$

3. **Path Independence:** For each $x, x', x'' \in \mathbb{X}$:

$$T_{x'',x'} \circ T_{x',x}(y) = T_{x'',x}(y), \quad \forall y \in \mathbb{Y}_x \quad (\text{PI})$$

Necessary Algebraic Properties

$$Y(x) = T_{x,x'}(Y(x')), \quad \forall x, x' \in \mathbb{X} \quad (\text{CC})$$

1. **Distribution Adaptedness:** For each $x, x' \in \mathbb{X}$,

$$(T_{x',x})_\# \mathbb{P}_{Y|X=x} = \mathbb{P}_{Y|X=x'} \quad (\text{DA})$$

2. **Identity:** For each $x \in \mathbb{X}$:

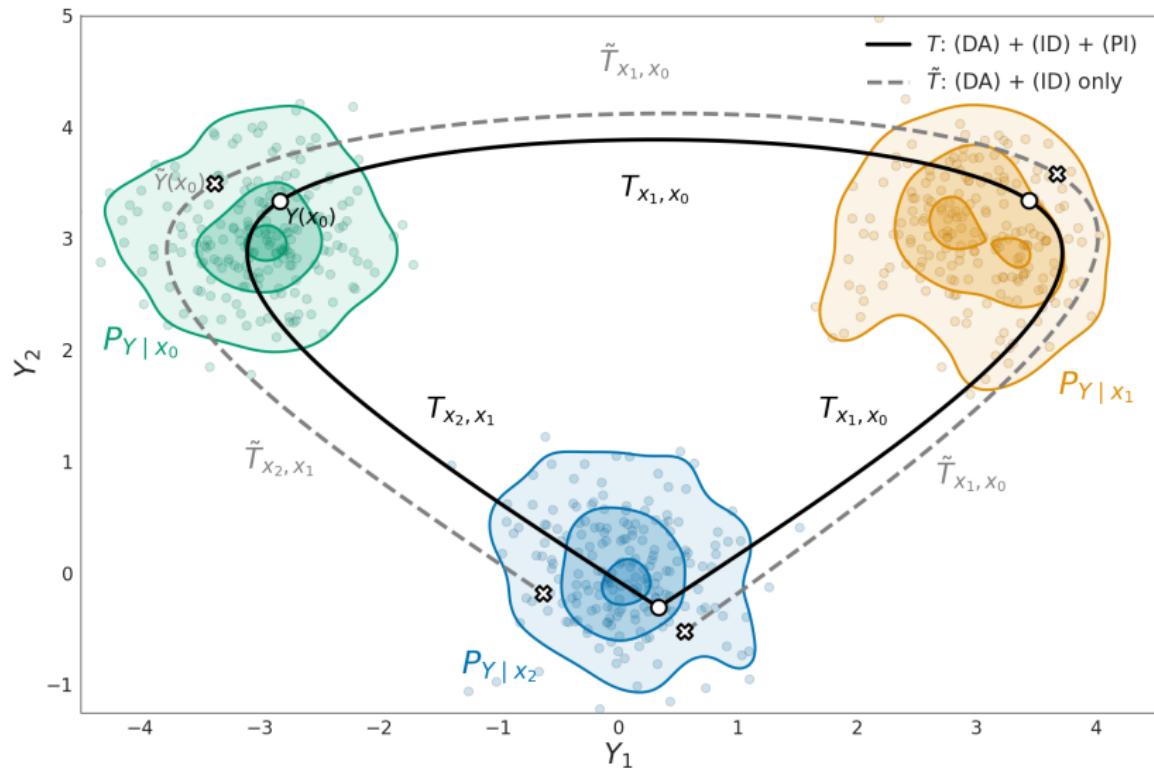
$$T_{x,x}(y) = y, \quad \forall y \in \mathbb{Y}_x \quad (\text{ID})$$

3. **Path Independence:** For each $x, x', x'' \in \mathbb{X}$:

$$T_{x'',x'} \circ T_{x',x}(y) = T_{x'',x}(y), \quad \forall y \in \mathbb{Y}_x \quad (\text{PI})$$

ID + PI = properties of a cocycle!

Importance of Path Independence



Sufficiency of Cocycle Properties for Admissible Transports

Are (DA) + (ID) + (PI) *sufficient* to enforce $Y(x) = T_{x,x'}(Y(x'))$? **not quite!**

But, they do guarantee the existence of *some* counterfactuals $\{\tilde{Y}(x)\}_{x \in \mathbb{X}}$:

$$\tilde{Y}(x) =_{\text{a.s.}} T_{x,x'}(\tilde{Y}(x')) \text{ for every } x, x' \in \mathbb{X}$$

Sufficiency of Cocycle Properties for Admissible Transports

Are (DA) + (ID) + (PI) *sufficient* to enforce $Y(x) = T_{x,x'}(Y(x'))$? **not quite!**

But, they do guarantee the existence of *some* counterfactuals $\{\tilde{Y}(x)\}_{x \in \mathbb{X}}$:

$$\tilde{Y}(x) =_{\text{a.s.}} T_{x,x'}(\tilde{Y}(x')) \text{ for every } x, x' \in \mathbb{X}$$

...how to enforce (ID) + (PI) + (DA) uniquely?

Structure of Counterfactual Cocycles

Theorem 1 (Cocycle Factorization)

Every cocycle T satisfying (ID), (PI), and (DA) w.r.t. $\mathbb{P}_{Y|X}$ can be represented as:

$$T_{x,x'} = f_x \circ f_{x'}^+, \quad \mathbb{P}_{Y|X=x'}\text{-a.s.}$$

Where $f_x : \mathbb{Y}_0 \rightarrow \mathbb{Y}$ is injective with left-inverse f_x^+

Structure of Counterfactual Cocycles

Theorem 1 (Cocycle Factorization)

Every cocycle T satisfying (ID), (PI), and (DA) w.r.t. $\mathbb{P}_{Y|X}$ can be represented as:

$$T_{x,x'} = f_x \circ f_{x'}^+, \quad \mathbb{P}_{Y|X=x'}\text{-a.s.}$$

Where $f_x : \mathbb{Y}_0 \rightarrow \mathbb{Y}$ is injective with left-inverse f_x^+

Proof Sketch

- Choose arbitrary $x_0 \in \mathbb{X}$
- Define $f_x := T_{x,x_0}$, $f_x^+ = T_{x_0,x}$

Structure of Counterfactual Cocycles

Theorem 1 (Cocycle Factorization)

Every cocycle T satisfying (ID), (PI), and (DA) w.r.t. $\mathbb{P}_{Y|X}$ can be represented as:

$$T_{x,x'} = f_x \circ f_{x'}^+, \quad \mathbb{P}_{Y|X=x'}\text{-a.s.}$$

Where $f_x : \mathbb{Y}_0 \rightarrow \mathbb{Y}$ is injective with left-inverse f_x^+

Proof Sketch

- Choose arbitrary $x_0 \in \mathbb{X}$
- Define $f_x := T_{x,x_0}$, $f_x^+ = T_{x_0,x}$

$$f_x \circ f_{x'}^+ = \underbrace{T_{x,x_0} \circ T_{x_0,x}}_{\text{PI}} = T_{x,x'}$$

Structure of Counterfactual Cocycles

Theorem 1 (Cocycle Factorization)

Every cocycle T satisfying (ID), (PI), and (DA) w.r.t. $\mathbb{P}_{Y|X}$ can be represented as:

$$T_{x,x'} = f_x \circ f_{x'}^+, \quad \mathbb{P}_{Y|X=x'}\text{-a.s.}$$

Where $f_x : \mathbb{Y}_0 \rightarrow \mathbb{Y}$ is injective with left-inverse f_x^+

Proof Sketch

- Choose arbitrary $x_0 \in \mathbb{X}$
- Define $f_x := T_{x,x_0}$, $f_x^+ = T_{x_0,x}$

$$f_x \circ f_{x'}^+ = \underbrace{T_{x,x_0} \circ T_{x_0,x}}_{\text{PI}} = T_{x,x'}$$

Implication: Can construct classes of cocycles via parameterized bijections!

$$\mathcal{F} \subseteq \{ f : \mathbb{X} \times \mathbb{Y} \rightarrow \mathbb{Y} \mid f_x := f(x, \bullet) \text{ bijective } \forall x \in \mathbb{X} \}$$

Examples: Normalizing flows, invertible NNs...

Can we use OT maps to construct a cocycle?

- Set $f_x := T_{x,x_0}^{(OT)}$ and $f_x^+ := T_{x_0,x}^{(OT)}$
- Since OT maps are invertible: $f_x^+ \circ f_x = \text{id} \implies \text{Left inverse!}$

Can we use OT maps to construct a cocycle?

- Set $f_x := T_{x,x_0}^{(OT)}$ and $f_x^+ := T_{x_0,x}^{(OT)}$
- Since OT maps are invertible: $f_x^+ \circ f_x = \text{id} \implies \text{Left inverse!}$

Since OT maps already satisfy (DA) + (NI) wrt. $P_{Y|X} \dots$

$$\hat{T}_{x,x'} = T_{x,x_0}^{(OT)} \circ T_{x_0,x'}^{(OT)} \implies \{\hat{T}_{x,x'}\}_{x,x' \in \mathbb{X}} \text{ satisfies (DA) + (ID) + (PI) wrt } \mathbb{P}_{Y|X}$$

Can we use OT maps to construct a cocycle?

- Set $f_x := T_{x,x_0}^{(OT)}$ and $f_x^+ := T_{x_0,x}^{(OT)}$
- Since OT maps are invertible: $f_x^+ \circ f_x = \text{id} \implies \text{Left inverse!}$

Since OT maps already satisfy (DA) + (NI) wrt. $P_{Y|X} \dots$

$$\hat{T}_{x,x'} = T_{x,x_0}^{(OT)} \circ T_{x_0,x'}^{(OT)} \implies \{\hat{T}_{x,x'}\}_{x,x' \in \mathbb{X}} \text{ satisfies (DA) + (ID) + (PI) wrt } \mathbb{P}_{Y|X}$$

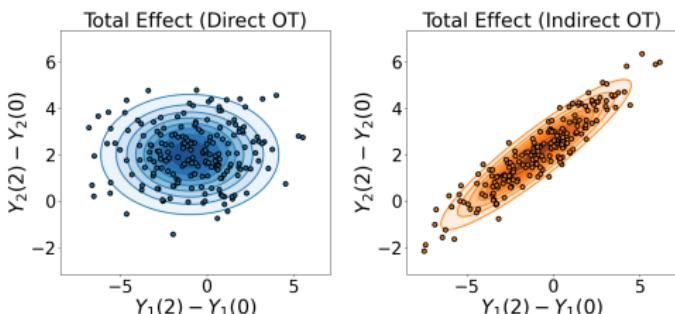
Problem: Construction depends on x_0 !

- Direct: $x_0 = 0$

$$T_{x,x'} = T_{x,0}^{(OT)} \circ T_{0,x'}^{(OT)}$$

- Indirect: $x_0 = 1$

$$T_{x,x'} = T_{x,1}^{(OT)} \circ T_{1,x'}^{(OT)}$$



Identifiability of Counterfactual Cocycles

Some Preliminaries

- If $f_x : \mathbb{Y} \rightarrow \mathbb{Y}$ is exactly invertible, $f_x \in \mathbb{G}$ (transformation group on \mathbb{Y})

Identifiability of Counterfactual Cocycles

Some Preliminaries

- If $f_x : \mathbb{Y} \rightarrow \mathbb{Y}$ is exactly invertible, $f_x \in \mathbb{G}$ (transformation group on \mathbb{Y})
- Transformation group \mathbb{G} means $f_1, f_2 \in \mathbb{G}$ then $f_1 \circ f_2 \in \mathbb{G}$ and $f_1^{-1}, f_2^{-1} \in \mathbb{G}$

Identifiability of Counterfactual Cocycles

Some Preliminaries

- If $f_x : \mathbb{Y} \rightarrow \mathbb{Y}$ is exactly invertible, $f_x \in \mathbb{G}$ (transformation group on \mathbb{Y})
- Transformation group \mathbb{G} means $f_1, f_2 \in \mathbb{G}$ then $f_1 \circ f_2 \in \mathbb{G}$ and $f_1^{-1}, f_2^{-1} \in \mathbb{G}$
- We call $f : (x, y) \mapsto f_x(y)$ a \mathbb{G} -valued coboundary map.

Identifiability of Counterfactual Cocycles

Some Preliminaries

- If $f_x : \mathbb{Y} \rightarrow \mathbb{Y}$ is exactly invertible, $f_x \in \mathbb{G}$ (transformation group on \mathbb{Y})
- Transformation group \mathbb{G} means $f_1, f_2 \in \mathbb{G}$ then $f_1 \circ f_2 \in \mathbb{G}$ and $f_1^{-1}, f_2^{-1} \in \mathbb{G}$
- We call $f : (x, y) \mapsto f_x(y)$ a \mathbb{G} -valued coboundary map.
- $\mathcal{F}_{\mathbb{G}}$ is the set of all \mathbb{G} -valued coboundary maps

$$\{f : \mathbb{X} \times \mathbb{Y} \rightarrow \mathbb{Y} \mid f(x, \bullet) := f_x \in \mathbb{G}\}$$

Identifiability of Counterfactual Cocycles

Some Preliminaries

- If $f_x : \mathbb{Y} \rightarrow \mathbb{Y}$ is exactly invertible, $f_x \in \mathbb{G}$ (transformation group on \mathbb{Y})
- Transformation group \mathbb{G} means $f_1, f_2 \in \mathbb{G}$ then $f_1 \circ f_2 \in \mathbb{G}$ and $f_1^{-1}, f_2^{-1} \in \mathbb{G}$
- We call $f : (x, y) \mapsto f_x(y)$ a \mathbb{G} -valued coboundary map.
- $\mathcal{F}_{\mathbb{G}}$ is the set of all \mathbb{G} -valued coboundary maps

$$\{f : \mathbb{X} \times \mathbb{Y} \rightarrow \mathbb{Y} \mid f(x, \bullet) := f_x \in \mathbb{G}\}$$

- We call $T_{x,x'}^{(f)} = f_x \circ f_{x'}^+$ a \mathbb{G} -valued cocycle

Identifiability of Counterfactual Cocycles

Some Preliminaries

- If $f_x : \mathbb{Y} \rightarrow \mathbb{Y}$ is exactly invertible, $f_x \in \mathbb{G}$ (transformation group on \mathbb{Y})
- Transformation group \mathbb{G} means $f_1, f_2 \in \mathbb{G}$ then $f_1 \circ f_2 \in \mathbb{G}$ and $f_1^{-1}, f_2^{-1} \in \mathbb{G}$
- We call $f : (x, y) \mapsto f_x(y)$ a \mathbb{G} -valued coboundary map.
- $\mathcal{F}_{\mathbb{G}}$ is the set of all \mathbb{G} -valued coboundary maps

$$\{f : \mathbb{X} \times \mathbb{Y} \rightarrow \mathbb{Y} \mid f(x, \bullet) := f_x \in \mathbb{G}\}$$

- We call $T_{x,x'}^{(f)} = f_x \circ f_{x'}^+$ a \mathbb{G} -valued cocycle
- Let $(\mathbb{P})|_{\mathbb{G}}$ be the set of transformations in \mathbb{G} that leave \mathbb{P} invariant.

Identifiability of Counterfactual Cocycles

Some Preliminaries

- If $f_x : \mathbb{Y} \rightarrow \mathbb{Y}$ is exactly invertible, $f_x \in \mathbb{G}$ (transformation group on \mathbb{Y})
- Transformation group \mathbb{G} means $f_1, f_2 \in \mathbb{G}$ then $f_1 \circ f_2 \in \mathbb{G}$ and $f_1^{-1}, f_2^{-1} \in \mathbb{G}$
- We call $f : (x, y) \mapsto f_x(y)$ a \mathbb{G} -valued coboundary map.
- $\mathcal{F}_{\mathbb{G}}$ is the set of all \mathbb{G} -valued coboundary maps

$$\{f : \mathbb{X} \times \mathbb{Y} \rightarrow \mathbb{Y} \mid f(x, \bullet) := f_x \in \mathbb{G}\}$$

- We call $T_{x,x'}^{(f)} = f_x \circ f_{x'}^+$ a \mathbb{G} -valued cocycle
- Let $(\mathbb{P})|_{\mathbb{G}}$ be the set of transformations in \mathbb{G} that leave \mathbb{P} invariant.

Theorem 2 (Identifiability of Counterfactual Cocycle)

Fix $x_0 \in \mathbb{X}$, and let T be a $\mathbb{P}_{Y|X}$ -adapted, \mathbb{G} -valued cocycle with coboundary map f .

$$T \text{ is identifiable in } \mathcal{F}_{\mathbb{G}} \Leftrightarrow (\mathbb{P}_{Y|X=x_0})|_{\mathbb{G}} \subseteq [id]_{\mathbb{P}_{Y|X=x_0}}$$

Smaller $\mathbb{G} \implies$ better for identifiability, worse for flexibility...

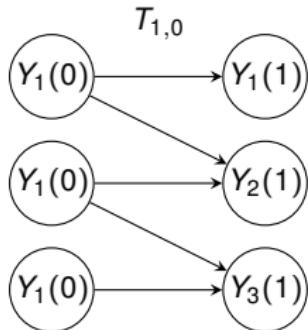
Practical Identifiability of Counterfactual Cocycles

Under known *causal ordering* of outcomes

$$Y_1 \prec Y_2 \prec \dots \prec Y_p$$

...natural to constrain $T_{x,x'}$ to be *lower triangular*:

$$f_x(y) = (f_{x,1}(y_1), \dots, f_{x,p}(y_1, \dots, y_p))$$



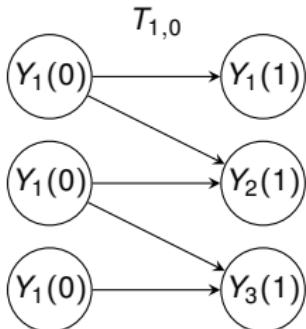
Practical Identifiability of Counterfactual Cocycles

Under known *causal ordering* of outcomes

$$Y_1 \prec Y_2 \prec \dots \prec Y_p$$

...natural to constrain $T_{x,x'}$ to be *lower triangular*:

$$f_x(y) = (f_{x,1}(y_1), \dots, f_{x,p}(y_1, \dots, y_p))$$



Enforcing conditional rank preservation leads to $f_x \in \mathbb{G}_{\text{TMI}}$!

Theorem 3 (Identifiability under TMI maps)

T is a $\mathbb{P}_{Y|X}$ -adapted, \mathbb{G}_{TMI} -valued cocycle $\implies T$ is identifiable in $\mathcal{F}_{\mathbb{G}_{\text{TMI}}}$.

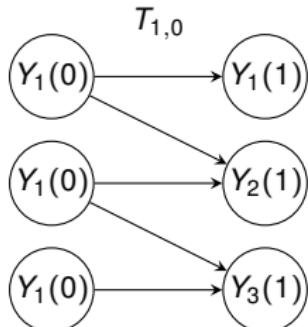
Practical Identifiability of Counterfactual Cocycles

Under known *causal ordering* of outcomes

$$Y_1 \prec Y_2 \prec \dots \prec Y_p$$

...natural to constrain $T_{x,x'}$ to be *lower triangular*:

$$f_x(y) = (f_{x,1}(y_1), \dots, f_{x,p}(y_1, \dots, y_p))$$



Enforcing conditional rank preservation leads to $f_x \in \mathbb{G}_{\text{TMI}}$!

Theorem 3 (Identifiability under TMI maps)

T is a $\mathbb{P}_{Y|X}$ -adapted, \mathbb{G}_{TMI} -valued cocycle $\implies T$ is identifiable in $\mathcal{F}_{\mathbb{G}_{\text{TMI}}}$.

- Can use autoregressive flows!
- Wait...isn't this just a flow-based BCM?

Model	Autoregressive transform $f_j(v_{<j}, \xi_j)$
NICE (Additive)	$\xi_j + \mu_j(v_{<j})$
MAF (Affine)	$\xi_j \mapsto \exp(\lambda_j(v_{<j})) \xi_j + \mu_j(v_{<j})$
NAF (INN)	$\xi_j \mapsto \sigma^{-1}(w(v_{<j}) \cdot \sigma(\sigma_j(v_{<j}) \xi_j + \mu_j(v_{<j})))$
NSF (Spline)	$\xi_j \mapsto v_j \mathbf{1}_{v_j \notin [-B, B]} + M_j(\xi_j; v_{<j}) \mathbf{1}_{v_j \in [-B, B]}$

Connection to SCMs

Theorem 4 (Cocycle Equivalence to Structural Model)

$\{Y(x)\}_{x \in \mathbb{X}}$ satisfies Exchangeability, Consistency and

$$Y(x) = T_{x,x'}(Y(x')) \quad \forall x, x' \in \mathbb{X}$$

if and only if $Y = f(X, \xi)$, $\xi \perp\!\!\!\perp X$.

Connection to SCMs

Theorem 4 (Cocycle Equivalence to Structural Model)

$\{Y(x)\}_{x \in \mathbb{X}}$ satisfies Exchangeability, Consistency and

$$Y(x) = T_{x,x'}(Y(x')) \quad \forall x, x' \in \mathbb{X}$$

if and only if $Y = f(X, \xi)$, $\xi \perp\!\!\!\perp X$.

Proof Sketch (\implies)

- Setting $f_x := T_{x,x_0}$ from before gives:

$$Y(x) = f_x \circ f_{x_0}^+(Y(x_0)) = f_x(Y(x_0)) =: f_x(\xi)$$

Connection to SCMs

Theorem 4 (Cocycle Equivalence to Structural Model)

$\{Y(x)\}_{x \in \mathbb{X}}$ satisfies Exchangeability, Consistency and

$$Y(x) = T_{x,x'}(Y(x')) \quad \forall x, x' \in \mathbb{X}$$

if and only if $Y = f(X, \xi)$, $\xi \perp\!\!\!\perp X$.

Proof Sketch (\implies)

- Setting $f_x := T_{x,x_0}$ from before gives:

$$Y(x) = f_x \circ f_{x_0}^+(Y(x_0)) = f_x(Y(x_0)) =: f_x(\xi)$$

- Exchangeability means $Y(x_0) \perp\!\!\!\perp X \implies \xi \perp\!\!\!\perp X$

Connection to SCMs

Theorem 4 (Cocycle Equivalence to Structural Model)

$\{Y(x)\}_{x \in \mathbb{X}}$ satisfies Exchangeability, Consistency and

$$Y(x) = T_{x,x'}(Y(x')) \quad \forall x, x' \in \mathbb{X}$$

if and only if $Y = f(X, \xi)$, $\xi \perp\!\!\!\perp X$.

Proof Sketch (\implies)

- Setting $f_x := T_{x,x_0}$ from before gives:

$$Y(x) = f_x \circ f_{x_0}^+(Y(x_0)) = f_x(Y(x_0)) =: f_x(\xi)$$

- Exchangeability means $Y(x_0) \perp\!\!\!\perp X \implies \xi \perp\!\!\!\perp X$
- Consistency means $Y(x) = f_x(\xi) \implies Y = f(X, \xi)$

Connection to SCMs

Theorem 4 (Cocycle Equivalence to Structural Model)

$\{Y(x)\}_{x \in \mathbb{X}}$ satisfies Exchangeability, Consistency and

$$Y(x) = T_{x,x'}(Y(x')) \quad \forall x, x' \in \mathbb{X}$$

if and only if $Y = f(X, \xi)$, $\xi \perp\!\!\!\perp X$.

Proof Sketch (\implies)

- Setting $f_x := T_{x,x_0}$ from before gives:

$$Y(x) = f_x \circ f_{x_0}^+(Y(x_0)) = f_x(Y(x_0)) =: f_x(\xi)$$

- Exchangeability means $Y(x_0) \perp\!\!\!\perp X \implies \xi \perp\!\!\!\perp X$
- Consistency means $Y(x) = f_x(\xi) \implies Y = f(X, \xi)$

Under TMI cocycle restriction, we recover Acyclic BCM $Y_j = f_j(Y_{<j}, X, \xi_j)$ again!

Our Idea: Center everything around the cocycle

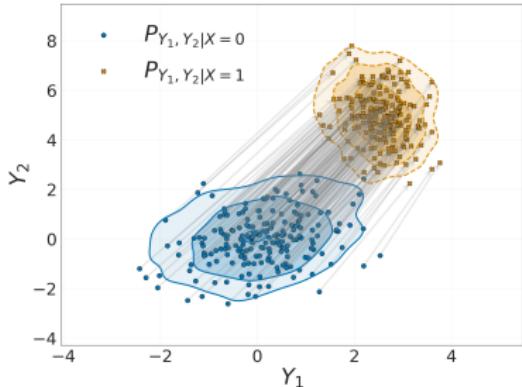
Aim: Do all estimation and inference without ever referencing or specifying \mathbb{P}_ξ

Our Idea: Center everything around the cocycle

Aim: Do all estimation and inference without ever referencing or specifying \mathbb{P}_{ξ}

1. Directly target $T_{x,x'}$ between conditionals:

$$\ell(T) = \mathbb{E}_{x,x'} D(\mathbb{P}_{Y|X=x}, (T_{x,x'})_{\#} \mathbb{P}_{Y|X=x'})^2$$

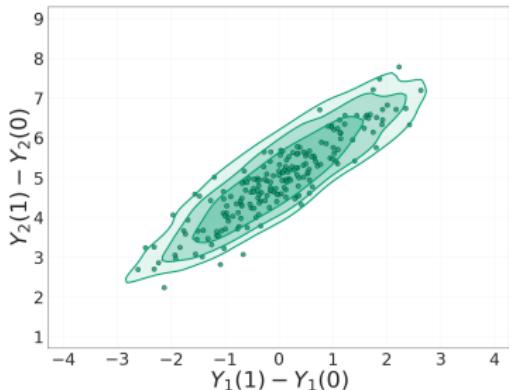


2. Use T to impute counterfactuals:

$$Y^{(i)}(x) = T_{x,X^{(i)}}(Y^{(i)})$$

... and empirically estimate quantities:

$$\widehat{THR} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(Y^{(i)}(1) \prec Y^{(i)}(0))$$



Benefits of Counterfactual Cocycle Modelling

Noise Invariance and Model Mis-Specification

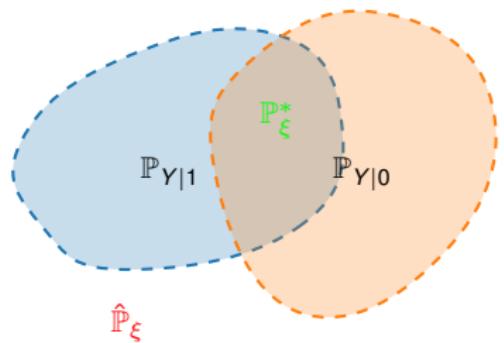
$$(\hat{f}_x)_\# : \hat{\mathbb{P}}_{\xi} \mapsto \mathbb{P}_{Y|X=x}, \quad (T_{x,x'})_\# = (f_x \circ f_{x'}^+)_\# : \mathbb{P}_{Y|X=x'} \mapsto \mathbb{P}_{Y|X=x}$$

Noise Invariance and Model Mis-Specification

$$(\hat{f}_x)_\# : \hat{\mathbb{P}}_\xi \mapsto \mathbb{P}_{Y|X=x}, \quad (T_{x,x'})_\# = (f_x \circ f_{x'}^+)_\# : \mathbb{P}_{Y|X=x'} \mapsto \mathbb{P}_{Y|X=x}$$

RCT case: $x \in \{0, 1\}$

- Consider class of functions \mathcal{F} for $(f_x)_{x \in \{0, 1\}}$
- Let $\mathcal{P}_1(\mathcal{F}), \mathcal{P}_0(\mathcal{F})$ = distributions reachable by pushing $\mathbb{P}_{Y|1}, \mathbb{P}_{Y|0}$ through each $f^+ \in \mathcal{F}^+$

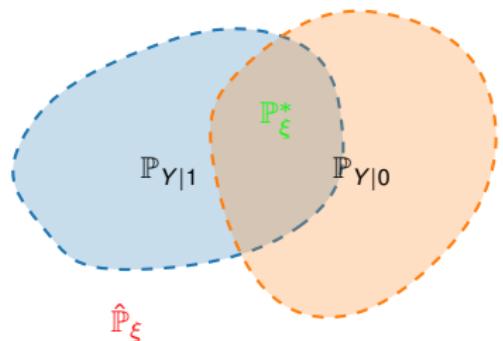


Noise Invariance and Model Mis-Specification

$$(\hat{f}_x)_\# : \hat{\mathbb{P}}_\xi \mapsto \mathbb{P}_{Y|X=x}, \quad (T_{x,x'})_\# = (f_x \circ f_{x'}^+)_\# : \mathbb{P}_{Y|X=x'} \mapsto \mathbb{P}_{Y|X=x}$$

RCT case: $x \in \{0, 1\}$

- Consider class of functions \mathcal{F} for $(f_x)_{x \in \{0, 1\}}$
- Let $\mathcal{P}_1(\mathcal{F}), \mathcal{P}_0(\mathcal{F})$ = distributions reachable by pushing $\mathbb{P}_{Y|1}, \mathbb{P}_{Y|0}$ through each $f^+ \in \mathcal{F}^+$



- $(\mathcal{F}, \hat{\mathbb{P}}_\xi)$ well-specified for (f, \mathbb{P}_ξ) if

$$\hat{\mathbb{P}}_\xi \in \mathcal{P}_1(\mathcal{F}) \cap \mathcal{P}_0(\mathcal{F})$$

- \mathcal{F} well specified for T if

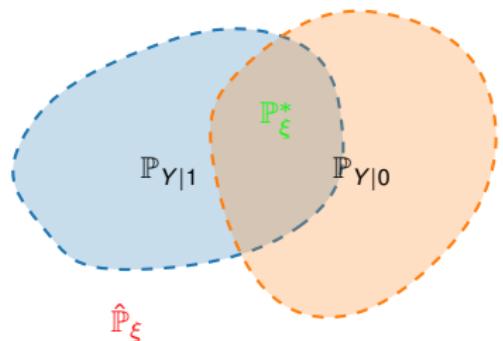
$$\exists \mathbb{P}_\xi^* \in \mathcal{P}_1(\mathcal{F}) \cap \mathcal{P}_0(\mathcal{F})$$

Noise Invariance and Model Mis-Specification

$$(\hat{f}_x)_\# : \hat{\mathbb{P}}_\xi \mapsto \mathbb{P}_{Y|X=x}, \quad (T_{x,x'})_\# = (f_x \circ f_{x'}^+)_\# : \mathbb{P}_{Y|X=x'} \mapsto \mathbb{P}_{Y|X=x}$$

RCT case: $x \in \{0, 1\}$

- Consider class of functions \mathcal{F} for $(f_x)_{x \in \{0, 1\}}$
- Let $\mathcal{P}_1(\mathcal{F}), \mathcal{P}_0(\mathcal{F})$ = distributions reachable by pushing $\mathbb{P}_{Y|1}, \mathbb{P}_{Y|0}$ through each $f^+ \in \mathcal{F}^+$



- $(\mathcal{F}, \hat{\mathbb{P}}_\xi)$ well-specified for (f, \mathbb{P}_ξ) if

$$\hat{\mathbb{P}}_\xi \in \mathcal{P}_1(\mathcal{F}) \cap \mathcal{P}_0(\mathcal{F})$$

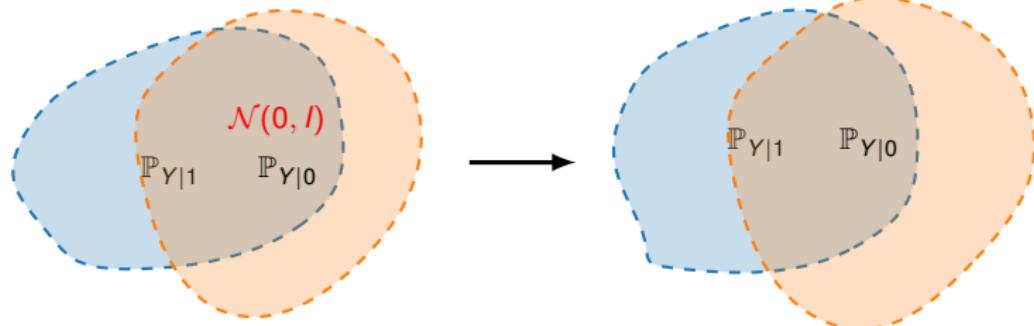
- \mathcal{F} well specified for T if

$$\exists \mathbb{P}_\xi^* \in \mathcal{P}_1(\mathcal{F}) \cap \mathcal{P}_0(\mathcal{F})$$

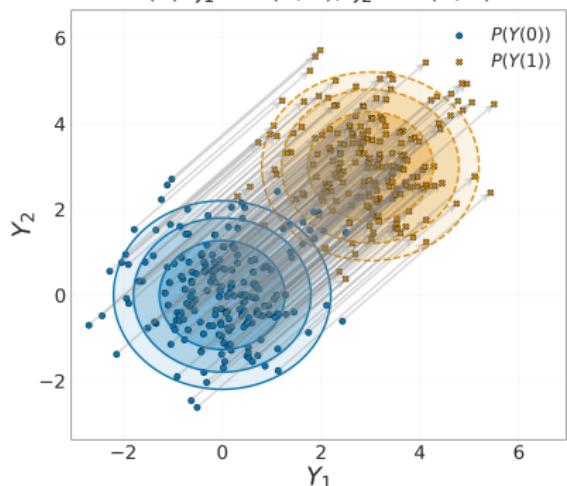
Invariance: Changing true \mathbb{P}_ξ doesn't change $\mathcal{P}_1(\mathcal{F}) \cap \mathcal{P}_0(\mathcal{F}) \neq \emptyset$!

$$\tilde{Y}(x) = f_x(\tilde{\xi}) \implies \tilde{Y}(x) = f_x \circ f_{x'}^+(\tilde{Y}(x'))$$

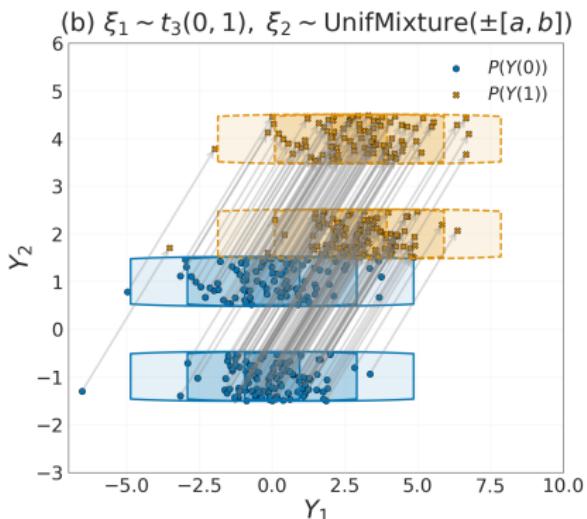
Invariance Example



(a) $\xi_1 \sim \mathcal{N}(0, 1)$, $\xi_2 \sim \mathcal{N}(0, 1)$



(b) $\xi_1 \sim t_3(0, 1)$, $\xi_2 \sim \text{UnifMixture}(\pm[a, b])$



Noise Invariance and Minimal Complexity

Can reparameterize SCM using any bijection $g \in \text{Aut}(\mathcal{E})$:

$$Y = f(X, \xi) = f(X, g \circ g^{-1}(\xi)) = f^{(g)}(X, \xi^{(g)})$$

Noise Invariance and Minimal Complexity

Can reparameterize SCM using any bijection $g \in \text{Aut}(\mathcal{E})$:

$$Y = f(X, \xi) = f(X, g \circ g^{-1}(\xi)) = f^{(g)}(X, \xi^{(g)})$$

Cocycle Invariant to Reparameterization:

$$T_{x,x'}^{(g)} = f_x^{(g)} \circ f_{x'}^{(g)} = f_x \circ \cancel{g \circ g^{-1}} \circ f_{x'}^+ = T_{x,x'}$$

Noise Invariance and Minimal Complexity

Can reparameterize SCM using any bijection $g \in \text{Aut}(\mathcal{E})$:

$$Y = f(X, \xi) = f(X, g \circ g^{-1}(\xi)) = f^{(g)}(X, \xi^{(g)})$$

Cocycle Invariant to Reparameterization:

$$T_{x,x'}^{(g)} = f_x^{(g)} \circ f_{x'}^{(g)} = f_x \circ \cancel{g \circ g^{-1}} \circ f_{x'}^+ = T_{x,x'}$$

Implication: Can use *any* member of $f^{(g)}$ to construct cocycle!

$$\implies \text{can choose } f^* \in \arg \min_f \text{Complexity}(f) \quad \text{s.t.} \quad f \in (f^{(g)})_{g \in \text{Aut}(\mathcal{E})}$$

Noise Invariance and Minimal Complexity

Can reparameterize SCM using any bijection $g \in \text{Aut}(\mathcal{E})$:

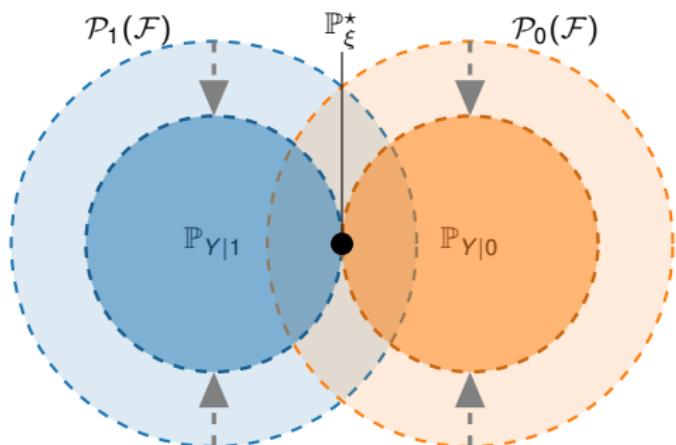
$$Y = f(X, \xi) = f(X, g \circ g^{-1}(\xi)) = f^{(g)}(X, \xi^{(g)})$$

Cocycle Invariant to Reparameterization:

$$T_{x,x'}^{(g)} = f_x^{(g)} \circ f_{x'}^{(g)} = f_x \circ \cancel{g \circ g^{-1}} \circ f_{x'}^+ = T_{x,x'}$$

Implication: Can use *any* member of $f^{(g)}$ to construct cocycle!

\implies can choose $f^* \in \arg \min_f \text{Complexity}(f)$ s.t. $f \in (f^{(g)})_{g \in \text{Aut}(\mathcal{E})}$



What is the optimal base distribution?

Recall $f_0, f_1 \in \mathbb{G} \implies T$ is \mathbb{G} -valued cocycle

Define smallest transformation group containing (f_0, f_1) : $\mathbb{G}_f := \langle f_0, f_1 \rangle$

What is the optimal base distribution?

Recall $f_0, f_1 \in \mathbb{G} \implies T$ is \mathbb{G} -valued cocycle

Define smallest transformation group containing $(f_0, f_1) : \mathbb{G}_f := \langle f_0, f_1 \rangle$

Size of \mathbb{G}_f : can depend on *parameterization* $(f_1^{(g)}, f_0^{(g)}) = (f^{(1)} \circ g, f^{(0)} \circ g)$

What is the optimal base distribution?

Recall $f_0, f_1 \in \mathbb{G} \implies T$ is \mathbb{G} -valued cocycle

Define smallest transformation group containing (f_0, f_1) : $\mathbb{G}_f := \langle f_0, f_1 \rangle$

Size of \mathbb{G}_f : can depend on parameterization $(f_1^{(g)}, f_0^{(g)}) = (f^{(1)} \circ g, f^{(0)} \circ g)$

Theorem 5 (Minimal Complexity Cocycle)

Take $(f_1^*, f_0^*) = (T_{1,0}, id)$. Then:

(i) $\mathbb{G}_{f^*} \subseteq \mathbb{G}_{f^* \circ g}$ and (ii) $\mathbb{G}_{f^*} \not\supseteq \mathbb{G}_{f^* \circ g}$ whenever $g \notin \mathbb{G}_{f^*}$

$\implies \mathbb{P}_{Y|0}$ is an optimal base distribution!

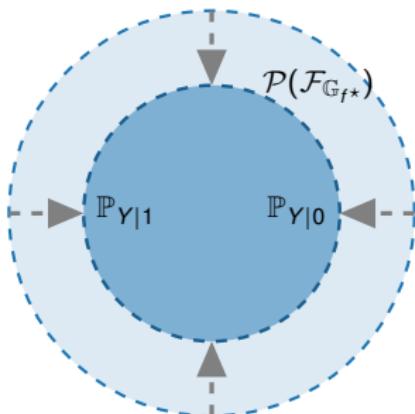
Example: $Y(x) = f_x(\xi) = x + \xi$

- Since $\mathbb{P}_\xi = \mathbb{P}_{Y|0}$:

$$f_x^* = f_x \in \mathbb{G} = (\mathbb{R}, +)$$

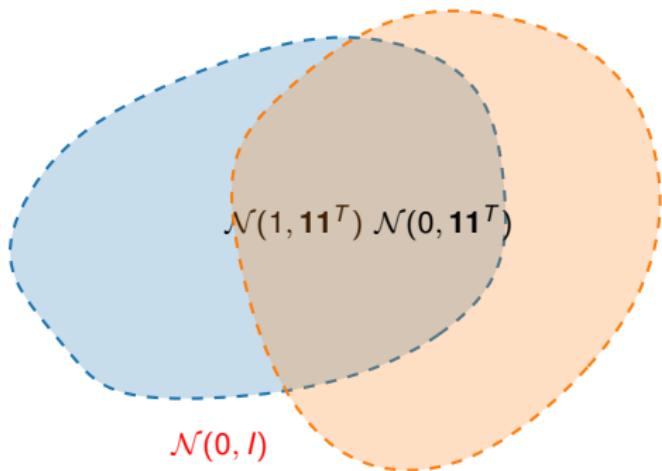
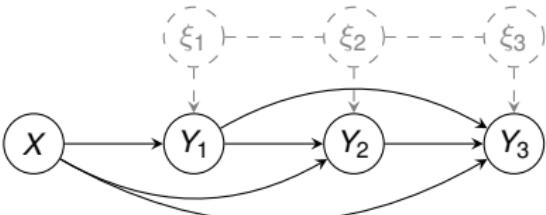
- Any \hat{f}_x not a shift of $\mathbb{P}_{Y|0}$ will need

$$\hat{f}_x \in \hat{\mathbb{G}} \not\subseteq (\mathbb{R}, +)$$



Robustness to Error Dependence

Counterfactual cocycles do *not* assume independent errors!



Example: Single latent cause

- $Y = \mathbf{1}X + \xi, \quad \xi \sim \mathcal{N}(0, \mathbf{1}\mathbf{1}^T)$
- If choosing $\hat{\mathbb{P}}_\xi = \mathcal{N}(0, I)$,
then even

$$\mathcal{F} = \{f : \mathbb{Y} \rightarrow \mathbb{Y} \text{ is bijective}\}$$

... is mis-specified for $f_0, f_1 \dots$

Cocycle Estimation

Estimating cocycles by minimising a distributional distance

Discrepancy in (DA):

$$\tilde{\ell}(T) = \mathbb{E}_{\textcolor{blue}{X}, \textcolor{red}{X'} \sim P_X} D(\mathbb{P}_{Y|X=\textcolor{blue}{X}}, (T_{\textcolor{blue}{X}, \textcolor{red}{X'}})_\# \mathbb{P}_{Y|X=\textcolor{red}{X'}})^2$$

Estimating cocycles by minimising a distributional distance

Discrepancy in (DA):

$$\tilde{\ell}(T) = \mathbb{E}_{\textcolor{blue}{X}, \textcolor{red}{X'} \sim P_X} D(\mathbb{P}_{Y|X=\textcolor{blue}{X}}, (T_{\textcolor{blue}{X}, \textcolor{red}{X'}})_\# \mathbb{P}_{Y|X=\textcolor{red}{X'}})^2$$

Issues:

1. Can't generally compute D in closed form since $\mathbb{P}_{Y|X}$ is unknown
2. No obvious empirical analogue $\tilde{\ell}_n(T) \rightarrow_p \tilde{\ell}(T)$.

Estimating cocycles by minimising a distributional distance

Discrepancy in (DA):

$$\tilde{\ell}(T) = \mathbb{E}_{X, X' \sim P_X} D(\mathbb{P}_{Y|X=X}, (T_{X, X'})_\# \mathbb{P}_{Y|X=X'})^2$$

Issues:

1. Can't generally compute D in closed form since $\mathbb{P}_{Y|X}$ is unknown
2. No obvious empirical analogue $\tilde{\ell}_n(T) \rightarrow_p \tilde{\ell}(T)$.

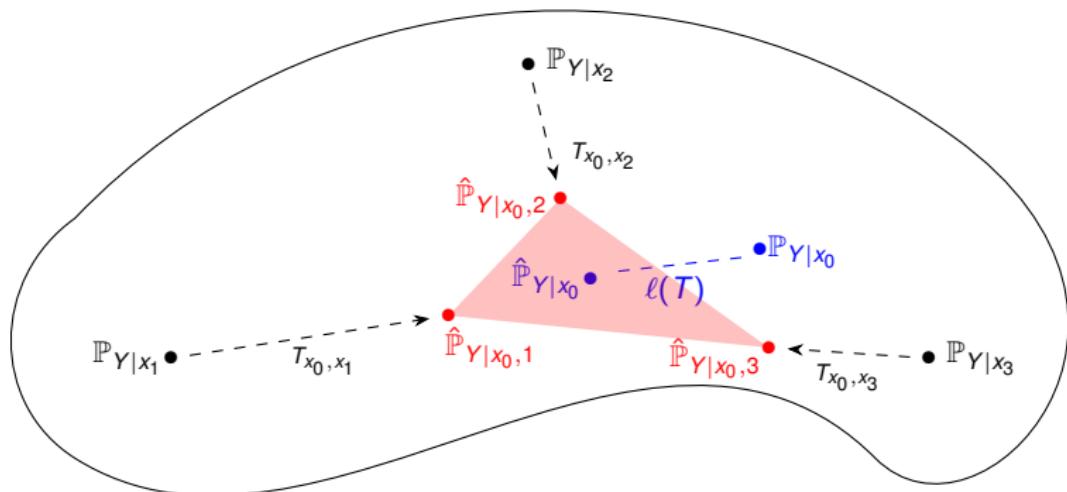
MMD Example: $\tilde{\ell}(T) = \mathbb{E} \| \mathbb{E}[\psi(Y)|X] - \mathbb{E}[\psi(T_{X,X}(Y')|X, X')] \|_{\mathcal{H}_k}^2$

... would need to learn conditional expectations $\mathbb{E}[\psi(Y)|X]$, $\mathbb{E}[\psi(T_{X,X}(Y')|X, X')]$...

A second attempt at the estimation objective

Idea: move expectation over X' inside the metric

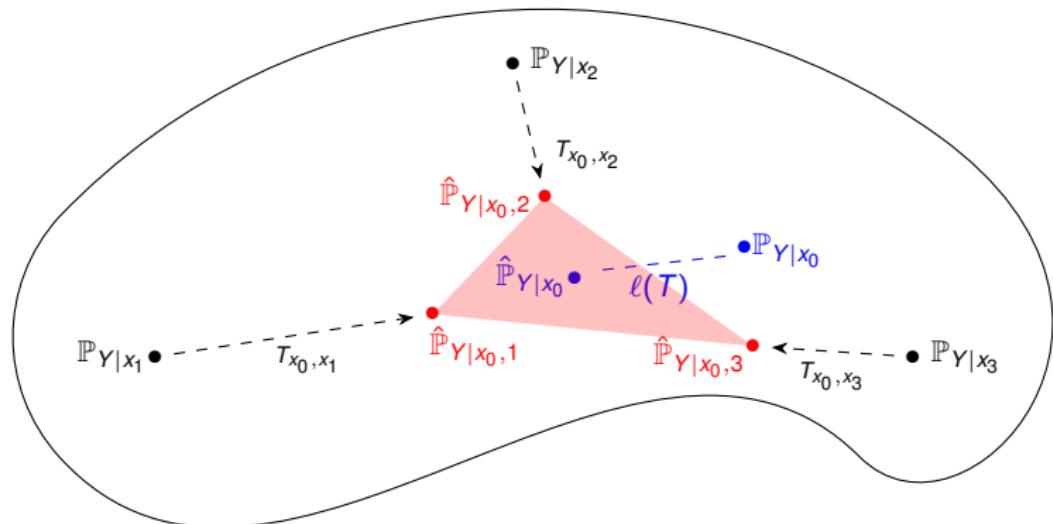
$$\ell(T) = \mathbb{E}_{X \sim P_X} D(\mathbb{P}_{Y|X=x}, \mathbb{E}_{X' \sim P_X} [((T_{X,X'}) \# \mathbb{P}_{Y|X=x'})]^2)$$



A second attempt at the estimation objective

Idea: move expectation over X' inside the metric

$$\ell(T) = \mathbb{E}_{X \sim P_X} D(\mathbb{P}_{Y|X=X}, \mathbb{E}_{X' \sim P_X} [(T_{X,X'}) \# \mathbb{P}_{Y|X=X'}])^2$$



- We prove $\arg \min_T \ell(T) = \arg \min_T \tilde{\ell}(T)$ in the well-specified setting!
- $\ell(T)$ can be estimated with simple empirical averages when $D = \text{MMD}$!

Cocycle Conditional Maximum Mean Discrepancy (CMMD)

CMMD: Take $D =$ as Maximum Mean Discrepancy:

$$\ell(T) = \mathbb{E} \left\| \psi(Y) - \mathbb{E}[\psi(T_{X,X'}(Y')|X)] \right\|_{\mathcal{H}}^2 + \text{constant}$$

Empirical analogue V-statistic :

$$\ell_n(T) = \frac{1}{n} \sum_i^n \left\| \psi(Y^{(i)}) - \frac{1}{n} \sum_j^n \psi(T_{X^{(i)}, X^{(j)}}(Y^{(j)})) \right\|_{\mathcal{H}}^2$$

Cocycle Conditional Maximum Mean Discrepancy (CMMD)

CMMD: Take $D =$ as Maximum Mean Discrepancy:

$$\ell(T) = \mathbb{E} \left\| \psi(Y) - \mathbb{E}[\psi(T_{X,X'}(Y')|X)] \right\|_{\mathcal{H}}^2 + \text{constant}$$

Empirical analogue V-statistic :

$$\ell_n(T) = \frac{1}{n} \sum_i^n \left\| \psi(Y^{(i)}) - \frac{1}{n} \sum_j^n \psi(T_{X^{(i)}, X^{(j)}}(Y^{(j)})) \right\|_{\mathcal{H}}^2$$

Theoretical properties

1. Under bounded, characteristic kernel k : $\ell_n(T) \rightarrow_{a.s.} \ell(T)$
2. Under continuous+compact+identifiable T : $\hat{T}_n \rightarrow_{a.s.} T$
3. Under regularity conditions on derivatives: \sqrt{n} consistency of parameterization θ
4. Unbiased minibatch-SGD with U-statistic

Cocycle Conditional Maximum Mean Discrepancy (CMMD)

CMMD: Take $D =$ as Maximum Mean Discrepancy:

$$\ell(T) = \mathbb{E} \left\| \psi(Y) - \mathbb{E}[\psi(T_{X,X'}(Y')|X)] \right\|_{\mathcal{H}}^2 + \text{constant}$$

Empirical analogue V-statistic :

$$\ell_n(T) = \frac{1}{n} \sum_i^n \left\| \psi(Y^{(i)}) - \frac{1}{n} \sum_j^n \psi(T_{X^{(i)}, X^{(j)}}(Y^{(j)})) \right\|_{\mathcal{H}}^2$$

Theoretical properties

1. Under bounded, characteristic kernel k : $\ell_n(T) \rightarrow_{a.s.} \ell(T)$
2. Under continuous+compact+identifiable T : $\hat{T}_n \rightarrow_{a.s.} T$
3. Under regularity conditions on derivatives: \sqrt{n} consistency of parameterization θ
4. Unbiased minibatch-SGD with U-statistic

In context of BCM $Y = f(X, \xi)$, consistency of \hat{T} does not depend on properties of ξ !

Extension to Larger Systems and Confounding

Current RCT setting implies:

$$Y = f(X, \xi), \quad X \perp\!\!\!\perp \xi, \quad f(X, \bullet) \text{ injective}$$

Limitations

- Injectivity will break if too many independent causes ξ_1, \dots, ξ_m
- Independence will break under confounding

Extension to Larger Systems and Confounding

Current RCT setting implies:

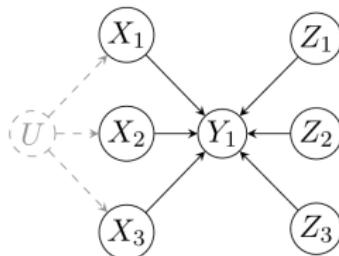
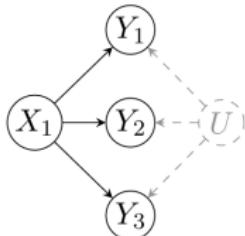
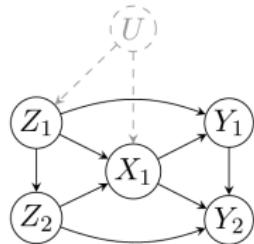
$$Y = f(X, \xi), \quad X \perp\!\!\!\perp \xi, \quad f(X, \bullet) \text{ injective}$$

Limitations

- Injectivity will break if too many independent causes ξ_1, \dots, ξ_m
- Independence will break under confounding

Extension: Assume set of measured causes Z of outcomes that satisfy:

$$Z \prec X \prec Y$$



Extension to Larger Systems and Confounding

Current RCT setting implies:

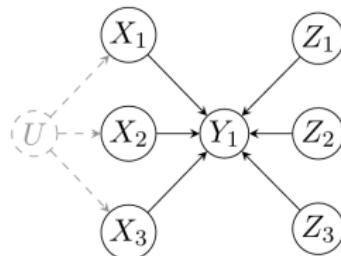
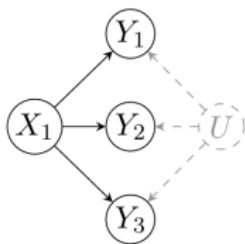
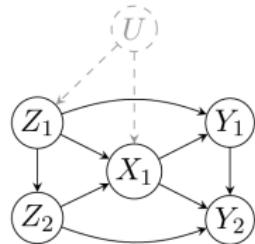
$$Y = f(X, \xi), \quad X \perp\!\!\!\perp \xi, \quad f(X, \bullet) \text{ injective}$$

Limitations

- Injectivity will break if too many independent causes ξ_1, \dots, ξ_m
- Independence will break under confounding

Extension: Assume set of measured causes Z of outcomes that satisfy:

$$Z \prec X \prec Y$$

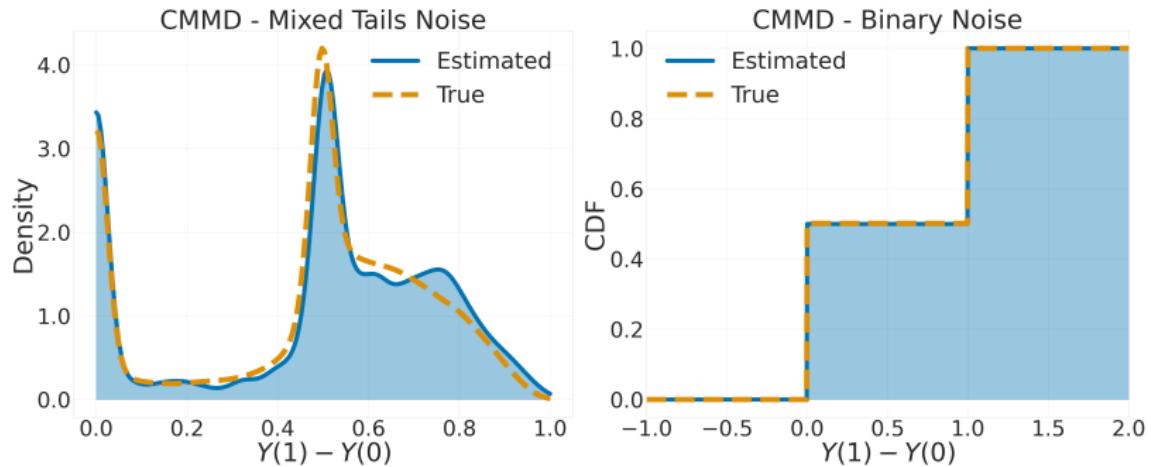


Idea: Do everything on counterfactuals $\{Y(x, z)\}_{(x, z) \in \mathbb{X} \times \mathbb{Z}}$ that relate to $Y(x)$ via nested consistency property: $Y(x) := Y(x, Z)$.

$$Y(x) = T_{(x, Z), (x', Z)}(Y(x'))$$

Experiments

Demonstration on the Toy Example



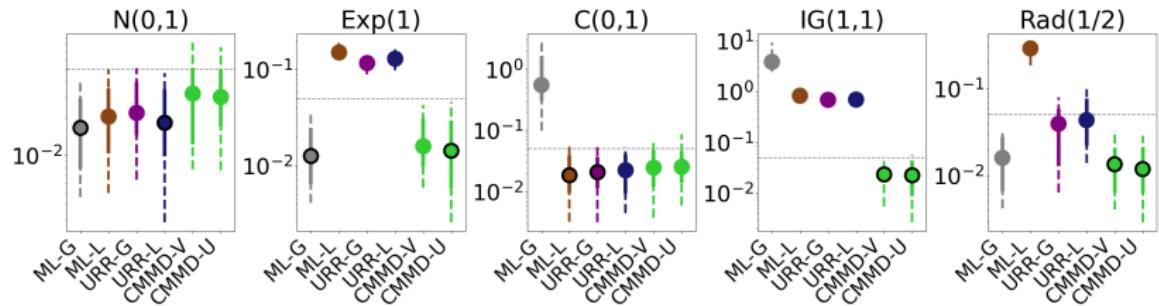
Noise Ablation vs flow-based SCMs

$$Y = X + \xi$$

Method	$\mathbb{P}_\xi = \mathcal{N}(0, 1)$	$\mathbb{P}_\xi = \mathcal{G}\mathcal{a}(1, 1)$	$\mathbb{P}_\xi = \mathcal{t}_1(0, 1)$	$\mathbb{P}_\xi = \mathcal{I}\mathcal{G}(1, 1)$	$\mathbb{P}_\xi = \text{Rad}(1/2)$
Interventional KS					
ML-G	0.015 ± 0.007	0.081 ± 0.041	0.121 ± 0.079	0.208 ± 0.162	0.405 ± 0.069
ML-L	0.055 ± 0.010	0.074 ± 0.036	0.110 ± 0.069	0.120 ± 0.090	0.415 ± 0.067
ML-T	0.018 ± 0.008	0.079 ± 0.040	0.029 ± 0.007	0.075 ± 0.031	0.412 ± 0.062
CMMMD-V	0.028 ± 0.009	0.027 ± 0.008	0.027 ± 0.009	0.027 ± 0.008	0.271 ± 0.031
CMMMD-U	0.010 ± 0.004	0.012 ± 0.005	0.008 ± 0.003	0.009 ± 0.004	0.268 ± 0.008
Counterfactual RMSE					
ML-G	0.035 ± 0.035	0.277 ± 0.118	113.667 ± 341.875	114.946 ± 147.841	0.326 ± 0.258
ML-L	0.036 ± 0.028	0.258 ± 0.073	97.745 ± 351.815	112.300 ± 165.014	0.480 ± 0.294
ML-T	0.033 ± 0.031	0.270 ± 0.097	0.044 ± 0.053	29.872 ± 41.628	0.391 ± 0.307
CMMMD-V	0.035 ± 0.026	0.020 ± 0.015	0.040 ± 0.031	0.028 ± 0.024	0.017 ± 0.019
CMMMD-U	0.040 ± 0.027	0.022 ± 0.016	0.033 ± 0.027	0.027 ± 0.023	0.014 ± 0.011
True Architecture Selection %					
ML-G	96%	14%	0%	2%	2%
ML-L	100%	2%	4%	0%	0%
ML-T	98%	8%	94%	0%	0%
CMMMD-V	100%	100%	100%	100%	98%
CMMMD-U	100%	100%	100%	100%	100%

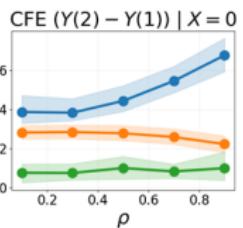
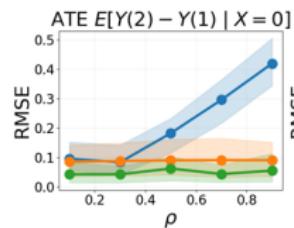
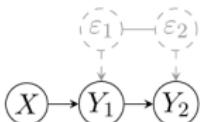
Noise Ablation vs flow-based SCMs

$$Y = X + \xi$$

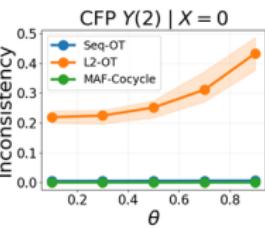
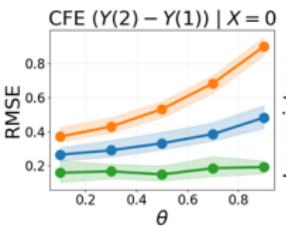
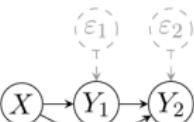


Confounding + Non-Additivity Ablation vs OT

Confounded Chain



Non-additive Triangle



Performance on SCM Benchmarks

Table 4: Mean \pm SD of KS_{int} and KS_{CF} on the *linear* SCMs.

Method	2var (lin)		triangle (lin)		fork (lin)		5chain (lin)	
	KS_{int}	KS_{CF}	KS_{int}	KS_{CF}	KS_{int}	KS_{CF}	KS_{int}	KS_{CF}
BGM	0.13 ± 0.07	0.19 ± 0.12	0.37 ± 0.05	0.06 ± 0.01	0.05 ± 0.07	0.61 ± 0.07	0.14 ± 0.05	0.06 ± 0.01
CausalNF	0.31 ± 0.11	0.24 ± 0.10	0.44 ± 0.05	0.06 ± 0.02	0.04 ± 0.02	0.66 ± 0.09	0.14 ± 0.02	0.06 ± 0.01
CAREFL	0.40 ± 0.04	0.15 ± 0.14	0.43 ± 0.05	0.06 ± 0.01	0.19 ± 0.01	0.58 ± 0.10	0.13 ± 0.02	0.06 ± 0.01
CocycleNF	0.03 ± 0.02	0.04 ± 0.04	0.23 ± 0.19	0.02 ± 0.01	0.02 ± 0.01	0.19 ± 0.23	0.02 ± 0.01	0.03 ± 0.01

Table 5: Mean \pm SD of KS_{int} and KS_{CF} on the *nonlinear* SCMs.

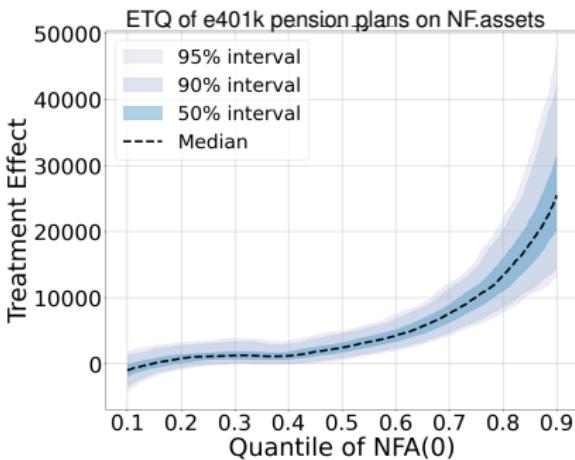
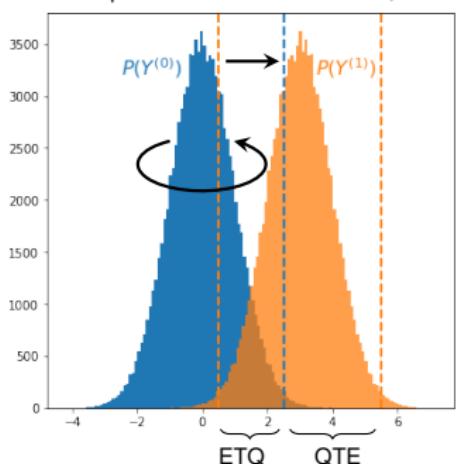
Method	2var (nonlin)		triangle (nonlin)		fork (nonlin)		5chain (nonlin)	
	KS_{int}	KS_{CF}	KS_{int}	KS_{CF}	KS_{int}	KS_{CF}	KS_{int}	KS_{CF}
BGM	0.12 ± 0.07	0.27 ± 0.13	0.41 ± 0.07	0.09 ± 0.05	0.04 ± 0.01	0.59 ± 0.08	0.07 ± 0.03	0.41 ± 0.12
CausalNF	0.30 ± 0.11	0.26 ± 0.08	0.47 ± 0.04	0.23 ± 0.08	0.07 ± 0.06	0.61 ± 0.09	0.06 ± 0.06	0.24 ± 0.16
CAREFL	0.44 ± 0.04	0.22 ± 0.15	0.47 ± 0.06	0.22 ± 0.09	0.19 ± 0.01	0.51 ± 0.21	0.17 ± 0.04	0.60 ± 0.25
CocycleNF	0.04 ± 0.02	0.13 ± 0.05	0.28 ± 0.13	0.20 ± 0.16	0.08 ± 0.05	0.20 ± 0.24	0.05 ± 0.02	0.20 ± 0.09

Counterfactual quantile effects on a real dataset

Quantile treatment effect: $\text{QTE}(\tau) = Q_{Y(1)}(\tau) - Q_{Y(0)}(\tau)$

Effect of Treatment on Quantile: $\text{ETQ}(\tau) = \mathbb{E}[Y^{(1)} - Y^{(0)} | Y^{(0)} = Q_{Y(0)}(\tau)]$

Example difference between ETQ,QTE



Summary

1. Counterfactual Cocycles as framework for Admissible counterfactual transports
2. Every counterfactual cocycle can be represented via left-invertible functions
3. Equivalent to injective SCMs, but cocycle is noise invariant + minimally complex
4. Robust cocycle estimator with consistency guarantees independent of true noise
5. Flexible parameterizations using flow-based toolkit
6. State-of-the-art performance on various benchmarks

Future Directions

1. Causal discovery?¹²
2. More general forms of confounding?
3. Non-iid settings?
4. Stochastic cocycles?
5. Efficiency theory?

¹²Xi, J., Dance, H., Orbanz, P. & Bloem-Reddy, B. '*Distinguishing Cause From Effect with Causal Velocity Models*' (ICML25).

Thank you!

Paper Link: [This Version](#) (soon to be on arXiv)

Github Repo: [hwdance/Cocycles](#)

To contact: hugh.dance.15@ucl.ac.uk, benbr@stat.ubc.ca



Appendix

