

Stochastic Causal Programming for Bounding Treatment Effects

<https://arxiv.org/pdf/2202.10806.pdf>

joint work with
Jakob Zeitler, David Watson, Matt Kusner,
Ricardo Silva, Niki Kilbertus

Kirtan Padh

HELMHOLTZ
TUM



KING'S
College
LONDON

Motivation

There was “a lot of correlation”

BRITISH

SMO

Membr

Professor of Medical Statistic

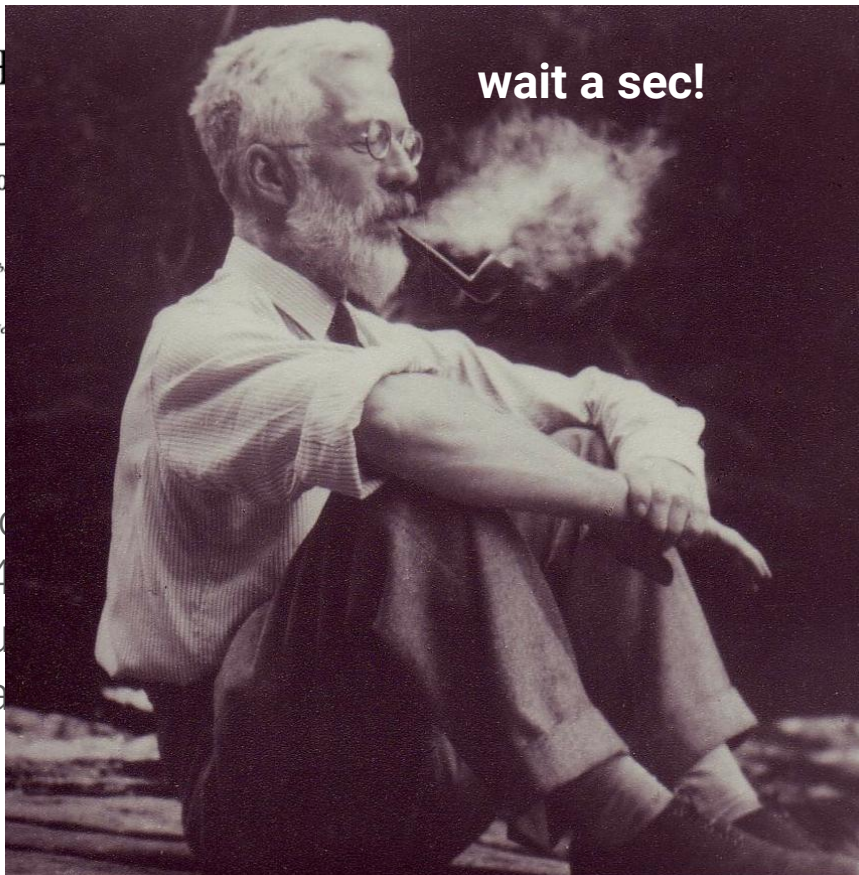
wait a sec!

**RELATIONSHIP BETWEEN HUMAN SMOKING
AND DEATH RATES**

LONG-TERM FOLLOW-UP STUDY OF 187,766 MEN

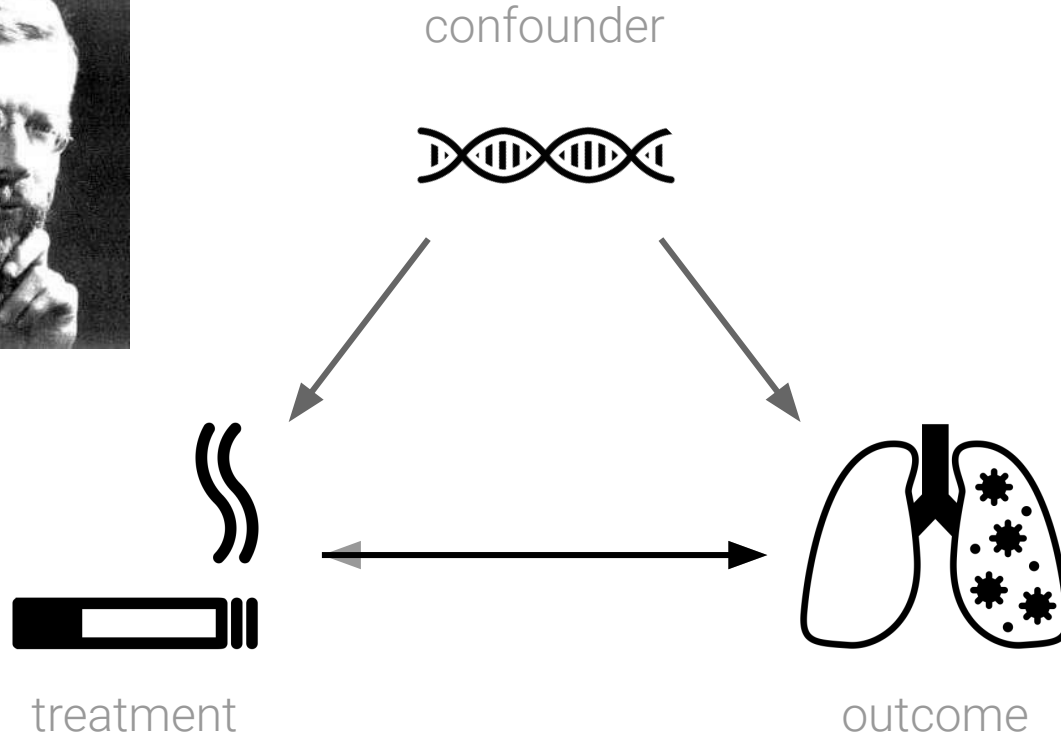
D.; Daniel Horn, Ph.D.

- 36
- 14
- su
- ca
- m

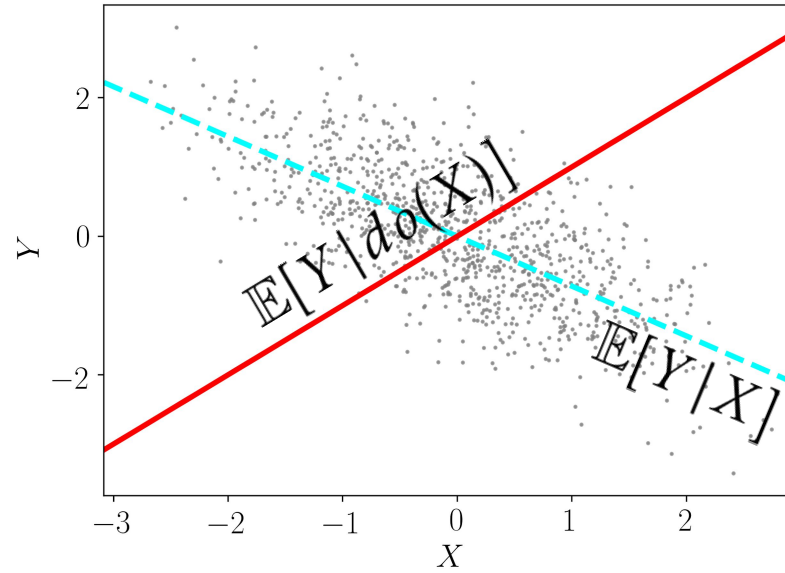


ers, 56 were heavy smokers
. 23.9% other cancer patients
st (all 36 who died of lung

Unobserved confounding

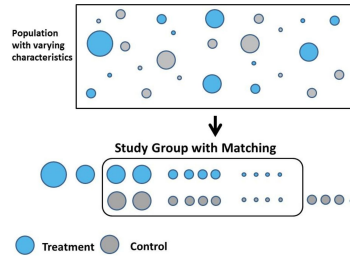


Naive ML approach fails



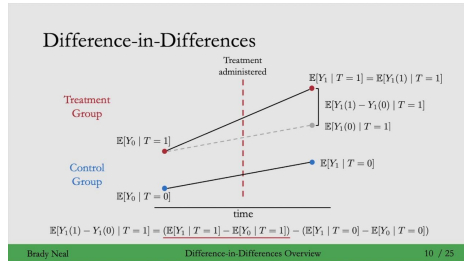
How do we estimate causal effects from observational data?

Some of the ways



Covariate matching

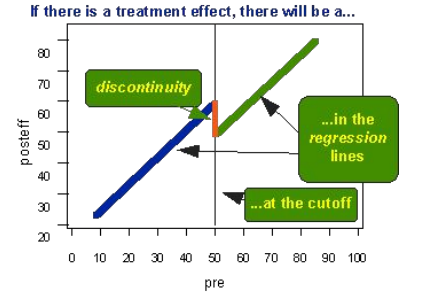
1



Difference-in-differences

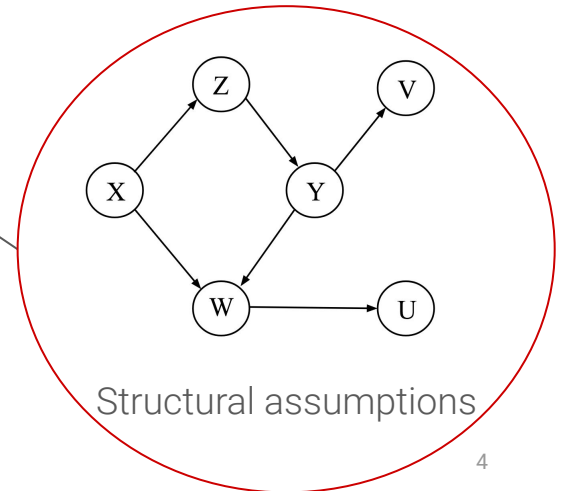
2

We are here!



Regression discontinuity

3



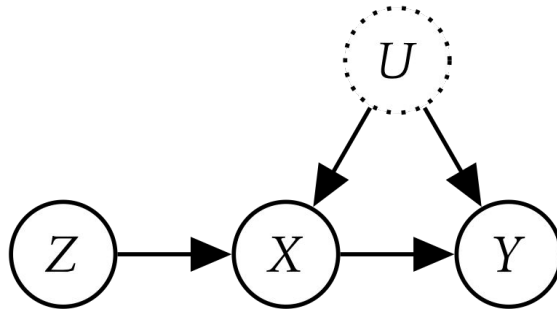
Structural assumptions

4

Structural assumptions

A known causal graph (DAG) with hidden confounding.

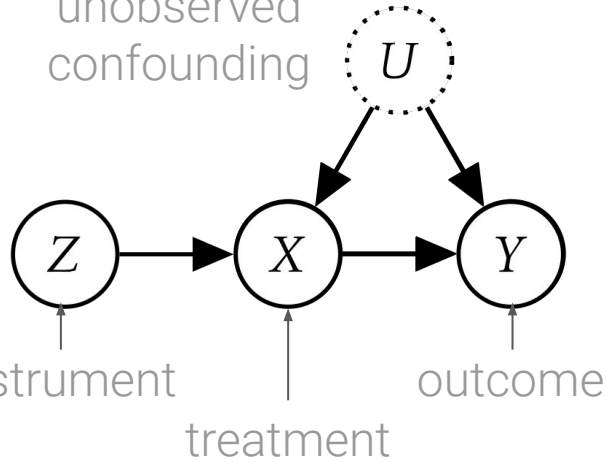
The method applies to general graphs



We focus on the **Instrument Variable (IV) setting**

IV: Identifiability (through additive noise)

unobserved
confounding



- (a) Z influences X $Z \not\perp\!\!\!\perp X$
- (b) Z is independent of U $Z \perp\!\!\!\perp U$
- (c) Z only influences Y via X $Z \perp\!\!\!\perp Y | \{X, U\}$

assume: $Y = f(X) + e_Y$ with $\mathbb{E}[e_Y] = 0$

$$\mathbb{E}[Y | z] = \mathbb{E}[f(X) + e_Y | z] = \mathbb{E}[f(X) | z] = \int f(x) p(x | z) dx$$

identifiable

unique under
mild conditions

identifiable

Even in the IV setting, conditions for identifiability are still (too?) strong

But

Assur

Point id

Addit

$$Y = f(X) + e_y$$

Mon

$$\Pr(x | z) \leq P$$

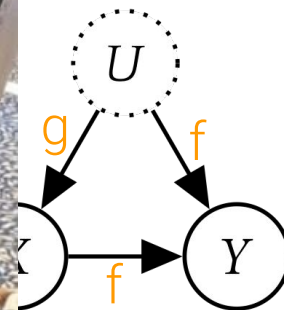
and

ability

IV setting

entifiability

nite number of
Gunsilius, 2019]



The story so far

- A zoo of point identification
- But less work on partial identification
 - Binary variables, I
 - Finite variables, g
 - Discrete variables
 - Scalar variables [K
 - High-dimensional

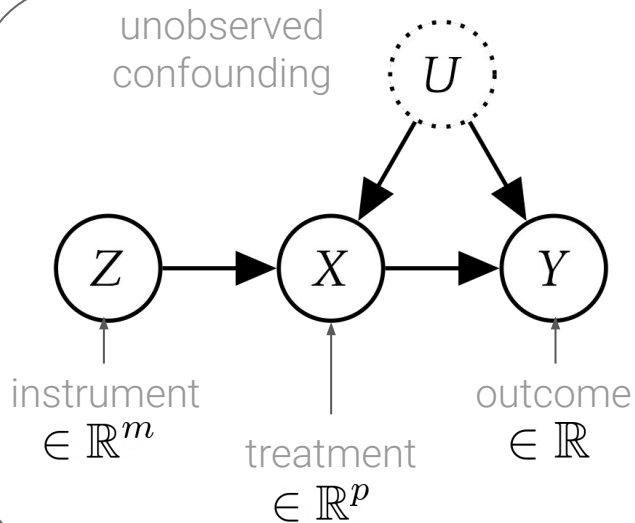


1930s.

21]

Problem formulation

General problem formulation



Assumptions

- (a) Z influences X $Z \not\perp\!\!\!\perp X$
- (b) Z is independent of U $Z \perp\!\!\!\perp U$
- (c) Z only influences Y via X $Z \perp\!\!\!\perp Y | \{X, U\}$

$$X = g(Z, U) \quad Y = f(X, U)$$

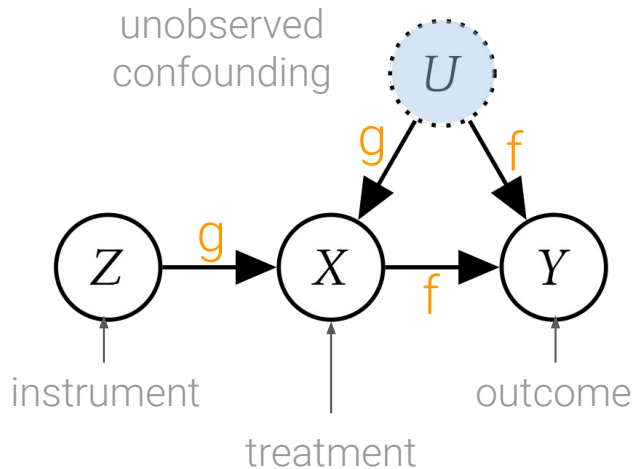
non-linear, non-additive

Goal - partial identification

For any x^* compute lower and upper bounds on the causal effect

$$\mathbb{E}[Y | do(x^*)]$$

General problem formulation as optimization



optimize over "all" distributions

$$X = g(Z, U) \quad Y = f(X, U)$$

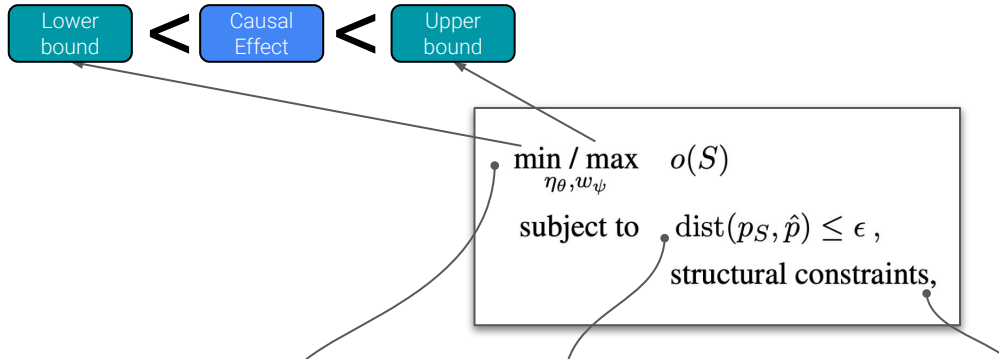
optimize over "all" functions

Detailed description: Two equations are shown. The first is $X = g(Z, U)$ and the second is $Y = f(X, U)$. A blue arrow points from the text 'optimize over "all" distributions' to the variable U in both equations. An orange arrow points from the text 'optimize over "all" functions' to the function symbols g and f in both equations.

Goal

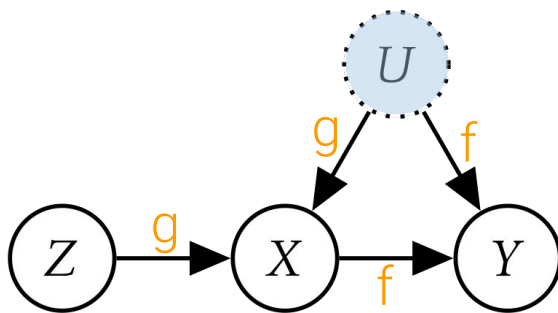
among all possible $\{g, f\}$ and distributions over U
that reproduce the observed densities $\{p(y | x, z), p(x | z)\}$,
estimate the min and max expected outcomes under intervention

A causal mathematical program



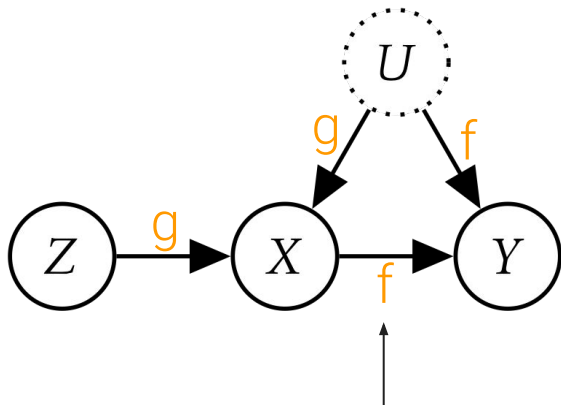
Cannot have no restrictions on f and g

- Without any restrictions on functions and distributions:
effect is not identifiable and average treatment effect bounds are vacuous
[Pearl, 1995; Bonet, 2001; Gunsilius 2018]
- Mild assumptions suffice for meaningful bounds:
 f and g have a finite number of discontinuities [Gunsilius, 2019]
- Rest of the talk: **operationalize the optimization**



Our practical approach

Response functions I [Balke & Pearl, 1994]



ultimately, we care about this functional relation

- Each value of U fixes a functional relation $X \rightarrow Y$
- Collect the set of all resulting functions $\{f_u\}$
- Identify values of u that result in the same f_u and assign a unique index r

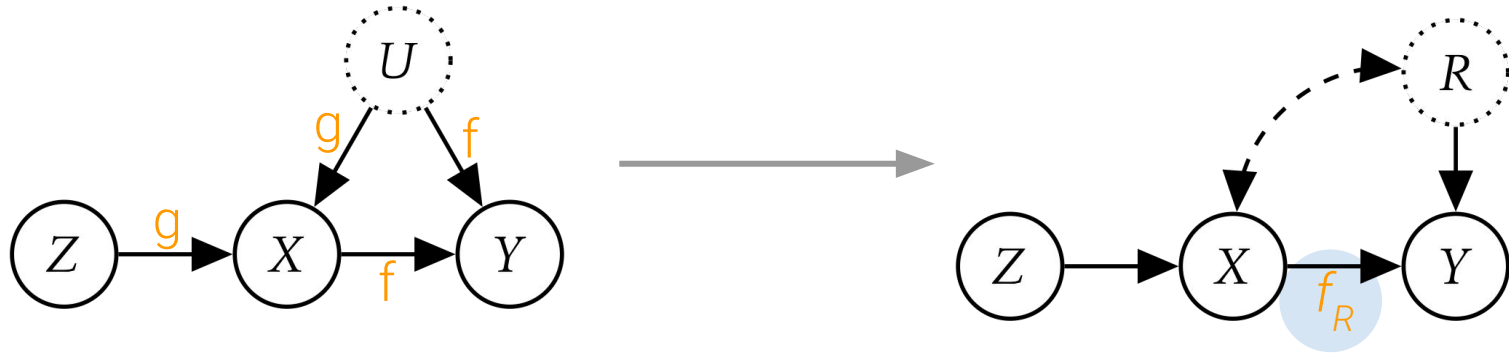
$$Y = f(X, U) = \lambda_1 X + \lambda_2 X U_1 + U_2$$

$$f(x, u) = \lambda_1 x + \lambda_2 x \quad \text{for } u_1 = 1, u_2 = 0$$

$$f_r(x) = (\lambda_1 + \lambda_2)x \quad \text{where } r \text{ is an alias for } (1, 0)$$

→ Instead of a potentially multivariate distribution over confounders U directly, we can think of a distribution R over functions $f: X \rightarrow Y$

Response functions II



choose convenient
function spaces

find convenient
representation of U from
which we can sample

find convenient representation of
distributions over response functions



Parameterizing response functions

We choose a simple parameterization

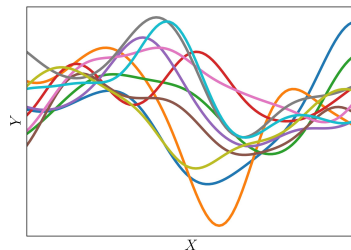
$$f_r(x) := f_{\theta_r}(x) \quad \text{for } \theta \in \Theta \subset \mathbb{R}^K$$

For simplicity, work with linear combination of (non-linear) basis functions:

$$f_{\theta}(x) = \sum_{k=1}^K \theta_k \psi_k(x) \quad \text{for basis functions } \{\psi_k : \mathbb{R}^p \rightarrow \mathbb{R}\}_{k \in [K]}$$



θ



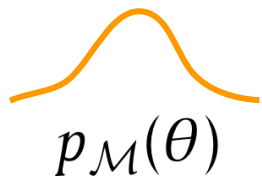
$f_{\theta} : X \rightarrow Y$

polynomials

neural networks

...

Parameterizing the distribution over θ



implies a causal model, and a distribution $\hat{p}_{\mathcal{M}}(x, y, z)$

Goal

Optimize over distributions $p_{\mathcal{M}}(\theta)$ such that

$\hat{p}_{\mathcal{M}}(x, y, z)$ is close to the observed distribution $p(x, y, z)$

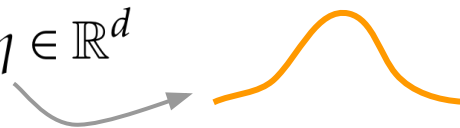
Ideally

low variance Monte-Carlo
gradient estimation

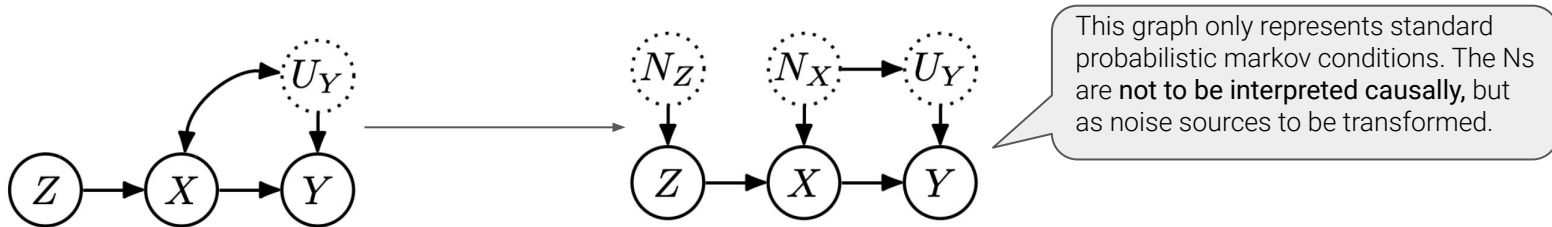
differentiable sampling

again, assume parametric form of $p_{\mathcal{M}}(\theta)$

$p_{\eta}(\theta)$ with $\eta \in \mathbb{R}^d$



The parametrization in practice



$$\theta \mid N_X \sim p_\eta(\cdot; \mu_{\eta_0}(N_X), \Sigma_{\eta_1}(N_X))$$

The distribution is defined up to mean and covariance functions.

Optimization Parameters: $\eta = (\eta_0, \eta_1)$

μ_{η_0} and Σ_{η_1} are small neural nets

A causal mathematical program

$$\begin{aligned} & \min / \max_{\eta_\theta, w_\psi} o(S) \\ & \text{subject to } \text{dist}(p_S, \hat{p}) \leq \epsilon, \\ & \text{structural constraints,} \end{aligned}$$

Objective: ATE [obj]

Optimising causal effect, e.g. ATE

Estimating $\mathbb{E}[Y|do(x^*)]$

Constraint: Data [c-data]

Matching the observed data distribution

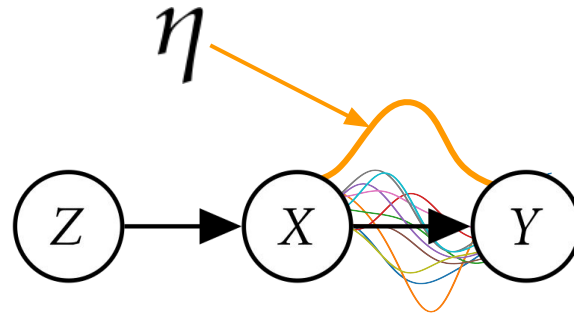
Matching $p(x | z)$ and $p(y | x, z)$

Constraint: Structure [c-struct]

Graphical assumptions

Baked into the model

Objective function



objective

$$\min_{\eta} / \max_{\eta} \mathbb{E}[Y | do(x^*)] = \min_{\eta} / \max_{\eta} \int f_{\theta}(x^*) p_{\eta}(\theta) d\theta$$
$$= \psi_Y(x^*)^{\top} \mathbb{E}_{N_X}[\mu_{\eta_0}(N_X)]$$

A causal mathematical program

$$\begin{array}{l} \min / \max_{\eta_\theta, w_\psi} o(S) \\ \text{subject to } \text{dist}(p_S, \hat{p}) \leq \epsilon, \\ \text{structural constraints,} \end{array}$$

Objective: ATE [obj]
Optimising causal effect, e.g. ATE

Constraint: Data [c-data]
Matching the observed data distribution

Constraint: Structure [c-struct]
Graphical assumptions

Estimating $\mathbb{E}[Y|do(x^*)]$

Matching $p(x | z)$ and $p(y | x, z)$

Baked into the model

Match $p(x | z)$

Identified from data and manually fixed once up front.
Implemented as an invertible **conditional normalizing flow**.

$$x = h_z(n)$$

Note: Given x_i, z_i , we can uniquely determine $n_i = h_{z_i}^{-1}(x_i)$

Match $p(y \mid x, z)$

Match the first two moments at a representative, finite set of points from $p(x, z)$

For the IV, the constraints are then in closed form.

$$\mathbb{E}[f_{\theta}(x_j) \mid x_j, z_j] = \psi(x_j)^{\top} \mu_{\eta_0}(n_j)$$

$$\mathbb{E}[f_{\theta}^2(x_j) \mid x_j, z_j] =$$

$$\psi(x_j)^{\top} \mu_{\eta_0}(n_j) (\Sigma_{\eta_1}(n_j) + \mu_{\eta_0}(n_j) \mu_{\eta_0}^{\top}(n_j)) \psi(x_j)$$

$$j \in \{1, 2, \dots, D\}$$

A random subsample from the data
(The 'stochastic' in stochastic causal programming)

A causal mathematical program

$$\min / \max_{\eta_\theta, w_\psi} o(S)$$

subject to $\text{dist}(p_S, \hat{p}) \leq \epsilon$,
structural constraints,

Objective: ATE [obj]

Optimising causal effect, e.g. ATE

Estimating $\mathbb{E}[Y|do(x^*)]$

Constraint: Data [c-data]

Matching the observed data distribution

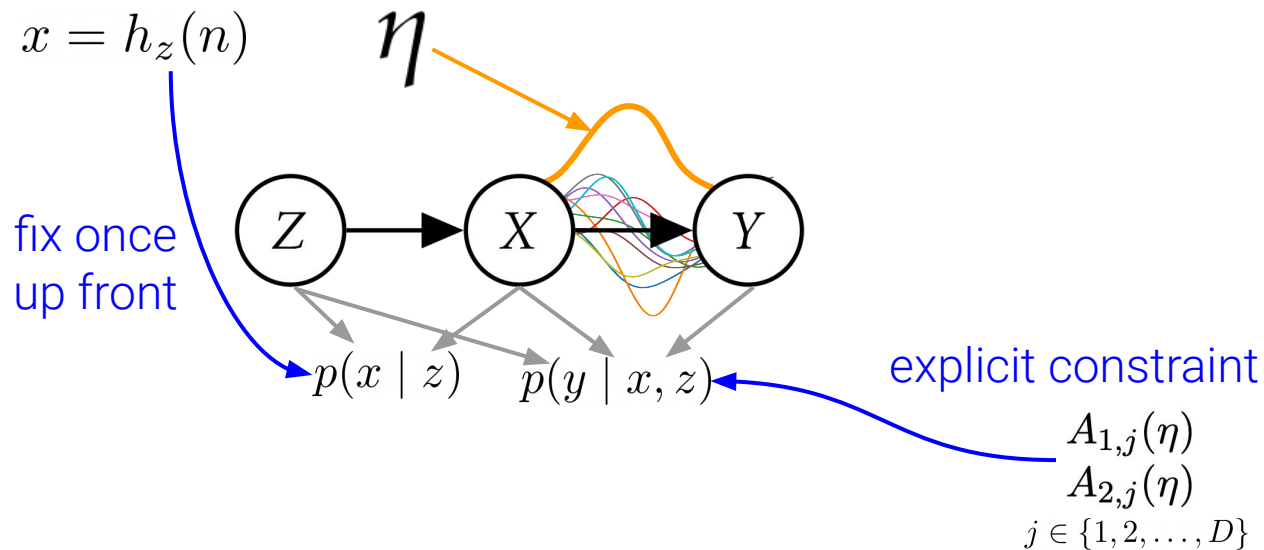
Matching $p(x | z)$ and $p(y | x, z)$

Constraint: Structure [c-struct]

Graphical assumptions

Baked into the model

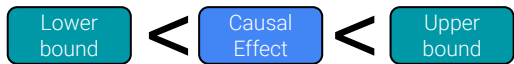
Intermediate overview



objective

$$\min_{\eta} / \max_{\eta} \mathbb{E}[Y | do(x^{\star})] :$$

The final optimization problem



Precompute once up front from data

The final values

$$\begin{aligned} & \min / \max_{\eta_\theta, w_\psi} o(S) \\ & \text{subject to } \text{dist}(p_S, \hat{p}) \leq \epsilon, \\ & \text{structural constraints,} \end{aligned}$$

Constraints enforced using the **augmented Lagrangian** (Nocedal and Wright, 2006)



Objective: [obj]

ATE

$$\psi_Y(x^*)^\top \mathbb{E}_{N_X} [\mu_{\eta_0}(N_X)]$$

Constraint: Data [c-data]

Matching the observed data distribution

$$x = h_z(n) \text{ (precompute } \{n_j\}_{j \in [D]})$$

$$\text{LHS}_{j,l} = \mathbb{E}[\phi_l(Y) \mid x_j, z_j]$$

$$\text{RHS}_{j,l}(\eta) = A_{j,l}$$

$$\begin{aligned} & \text{LHS}_{j,l} = \text{RHS}_{j,l}(\eta) \\ & \forall l \in [2] \text{ and } j \in [D] \end{aligned}$$

Constraint: Structure [c-struct]

Graphical assumptions

Baked into the model

Empirical results

Choices of response functions

$$f_{\theta}(x) = \sum_{k=1}^K \theta_k \psi_k(x) \text{ for basis functions } \{\psi_k : \mathbb{R}^p \rightarrow \mathbb{R}\}_{k \in [K]}$$

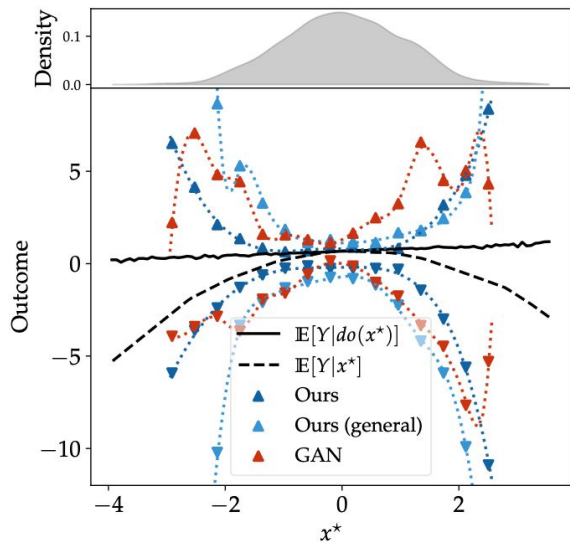
Polynomials

Neural network

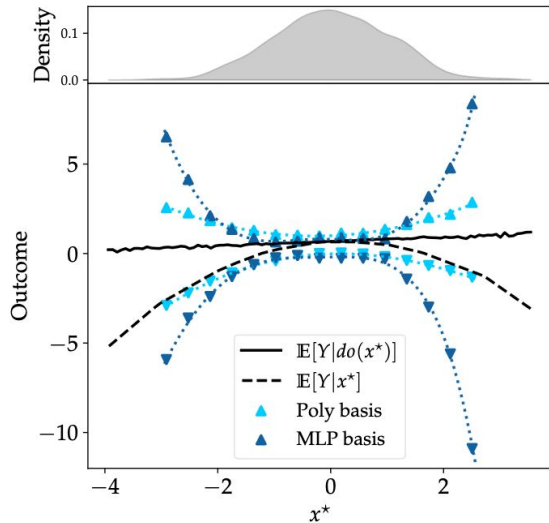
Train a small fully connected network on observed data $X \rightarrow Y$ and take activations of last hidden layer as basis functions.

For visualization: All interventions are along a single axis for multi-dimensional treatments

Linear 2-D treatment,
strong confounding

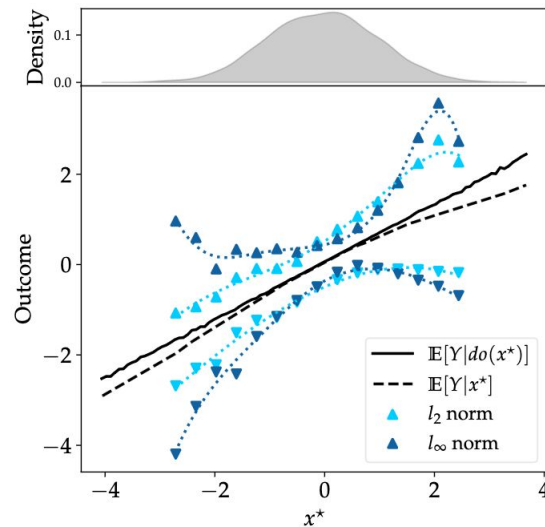


Ours seems more stable



Stronger assumptions, stronger inference

Linear 2-D treatment,
weak confounding



$$\sum_{j,l} (\text{LHS}_{j,l} - \text{RHS}_{j,l})^2 = 0$$

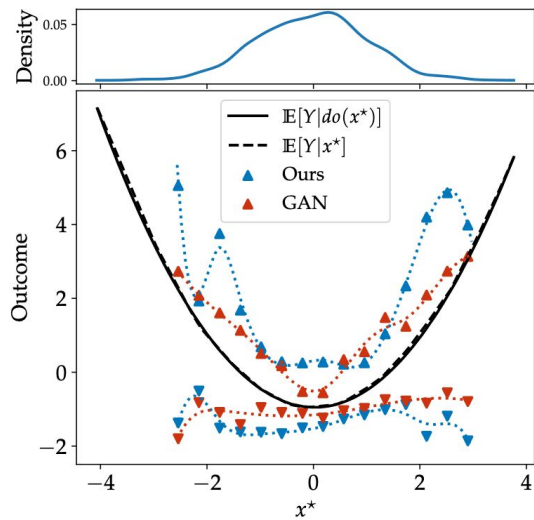
Single constraint [c-data]

Takeaways

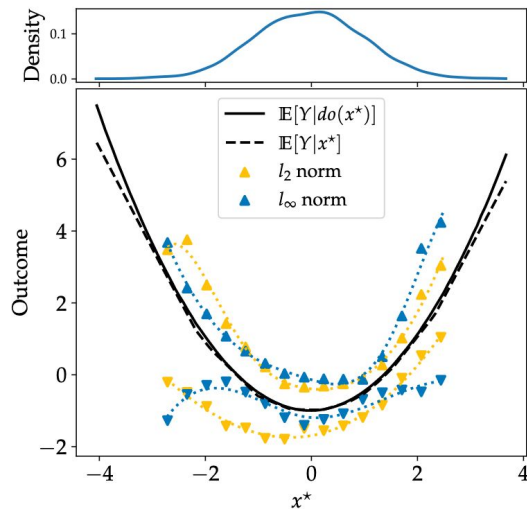
- Partial identification is hard for continuous high-dimensional variables.
- We were able to craft a framework that is
 - **flexible**
 - **needs minimal Monte-Carlo estimations** in the IV and leaky mediator settings
 - allows the user the choice of a **spectrum of assumption strength**

See the paper for extensions to more general settings.

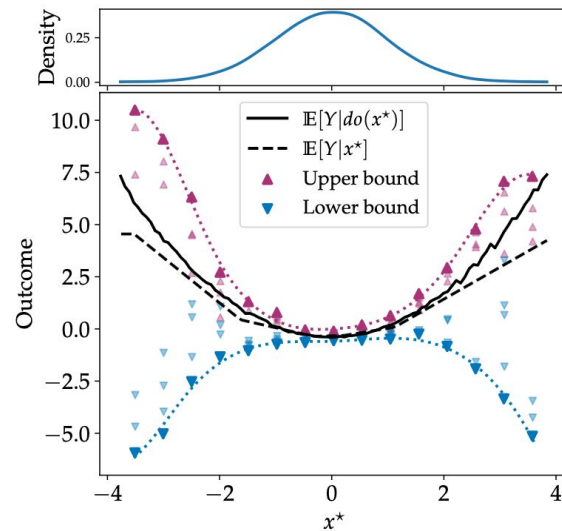
Quadratic 3-D treatment,
weak confounding



Quadratic 2-D treatment,
weak confounding



Quadratic scalar treatment,
weak confounding



Thank you