

A class of algorithms for general instrumental variable models

<https://arxiv.org/abs/2006.06366>
(NeurIPS 2020)

joint work with
Matt Kusner & **Ricardo Silva**



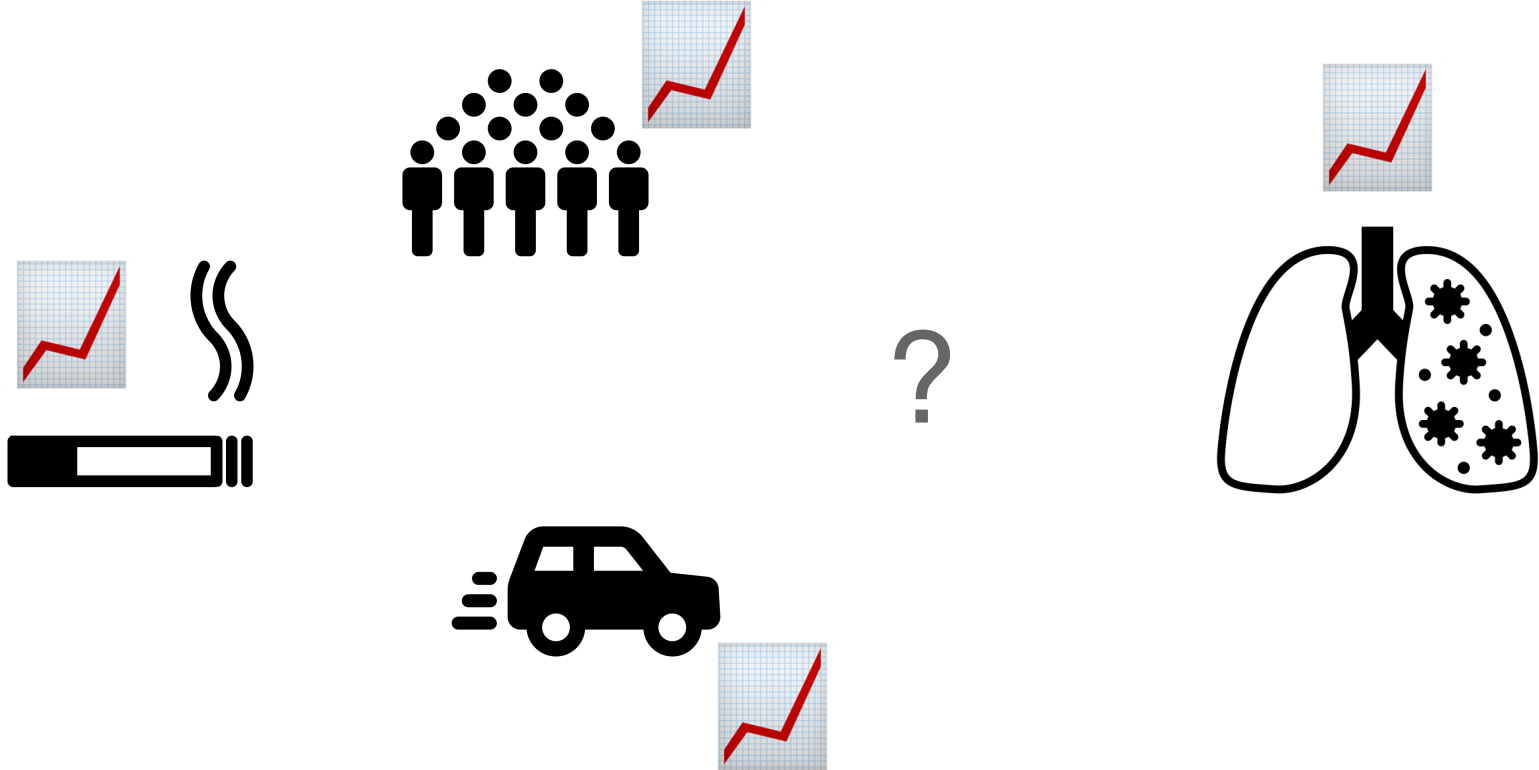
Niki Kilbertus

HELMHOLTZAI



Motivation

Let's start with a classic



There was “a lot of correlation”

BRITISH

SMO

Mem

Professor of Medical Statistic

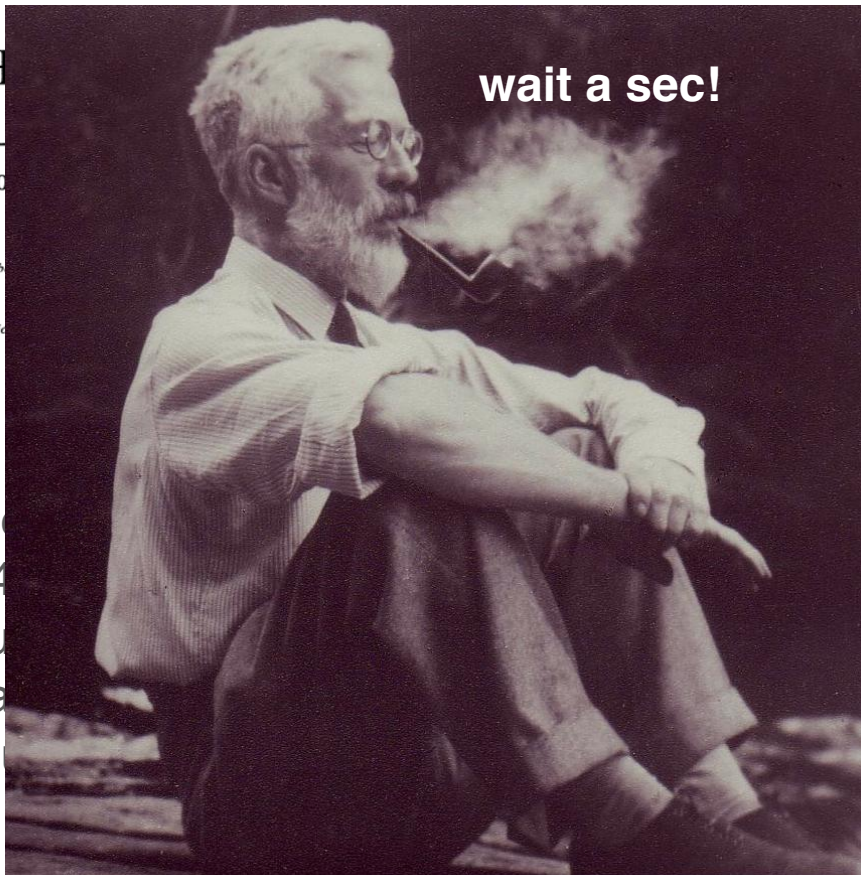
wait a sec!

RELATIONSHIP BETWEEN HUMAN SMOKING AND DEATH RATES

W-UP STUDY OF 187,766 MEN

D.; Daniel Horn, Ph.D.

- 36
 - 14
 - su
 - ca
 - m
- ers, 56 were heavy smokers
s. 23.9% other cancer patients
st (all 36 who died of lung



Unobserved confounding

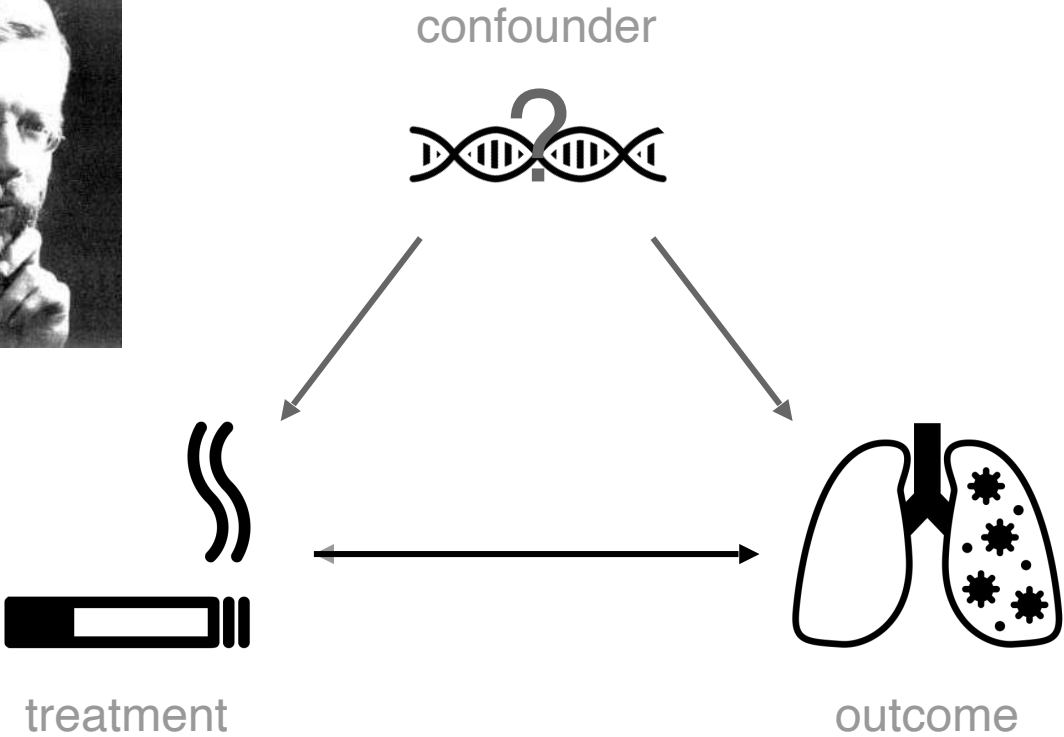
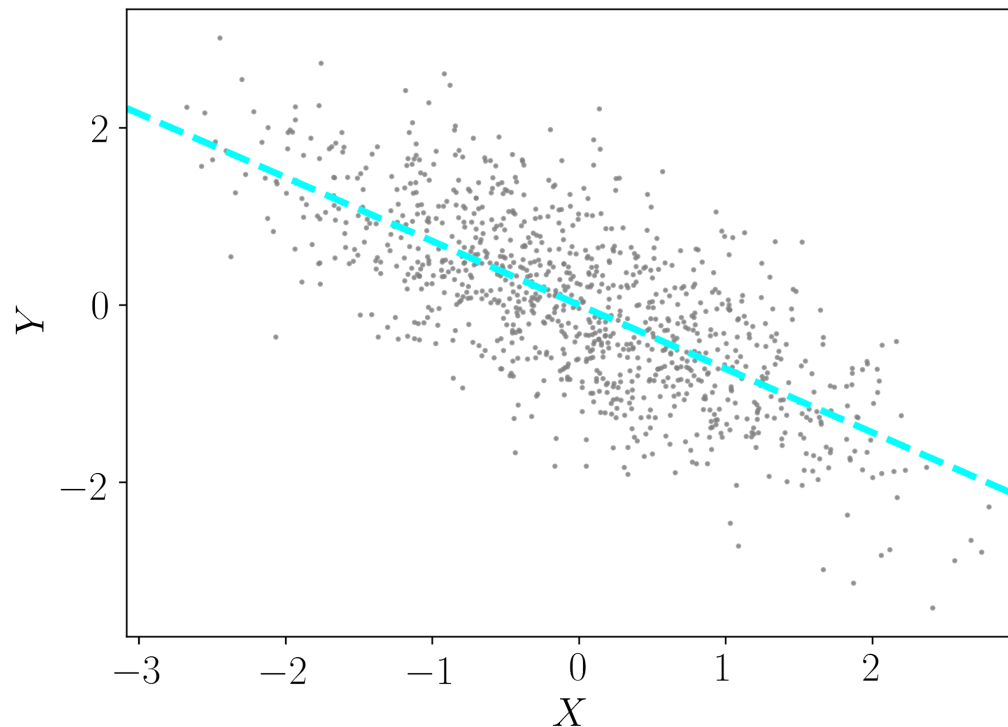


Image credit (from Noun Project):
priyanka, Andrew Nielsen, mungang kim

Introduction

Naive ML approach: standard regression



$$X \in \mathbb{R}, Y \in \mathbb{R}$$

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

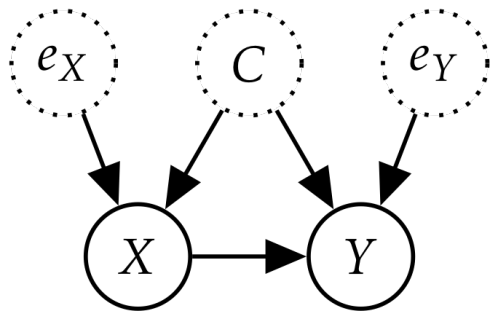
$$\mathbb{E}[e_Y | X] = 0$$

$$\mathbb{E}[Y - f(X) | X] = 0$$

$$\Rightarrow \mathbb{E}[Y | X] = f(X)$$

linear least squares:
$$f = \arg \min_{\hat{f}} \sum_i (\hat{f}(x_i) - y_i)^2$$

Naive ML approach failing

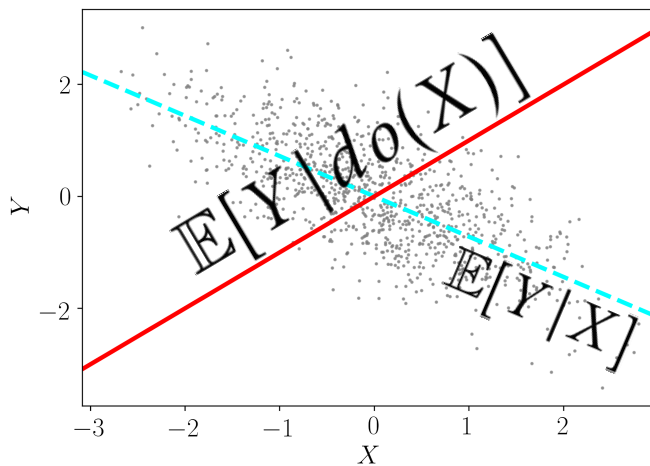


$$X = \alpha \cdot C + e_X$$

$$Y = X + \beta \cdot C + e_Y$$

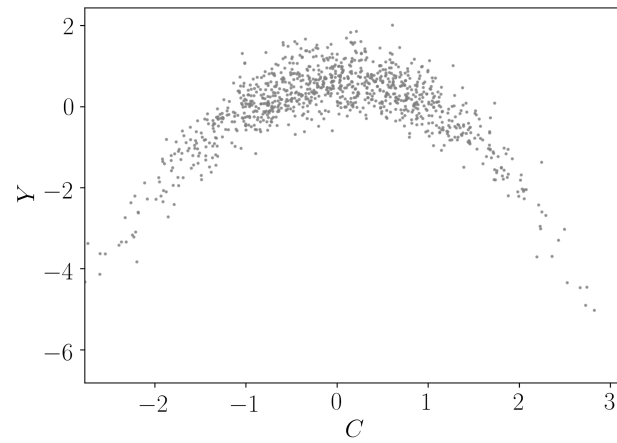
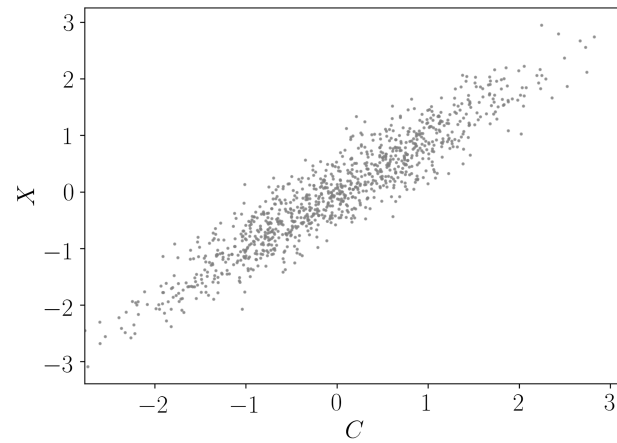
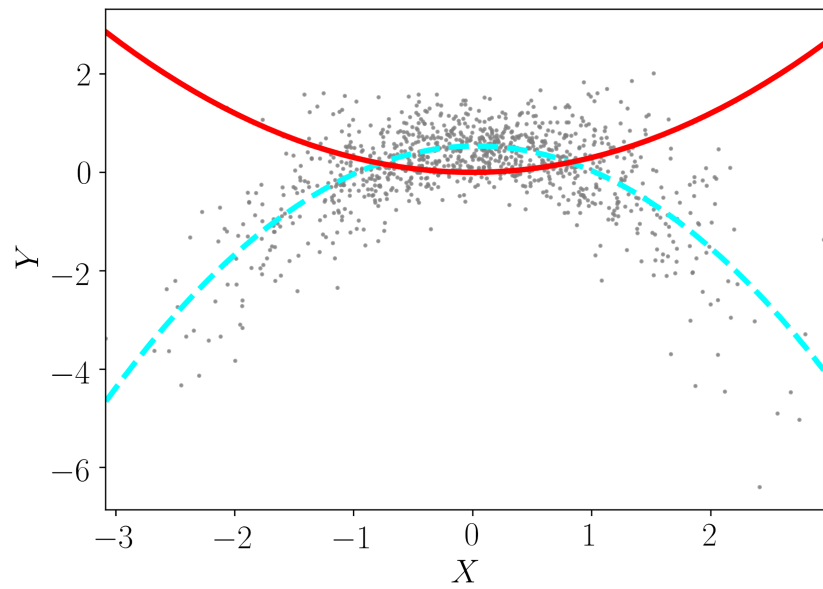


$$Y = f(X) + e_Y$$

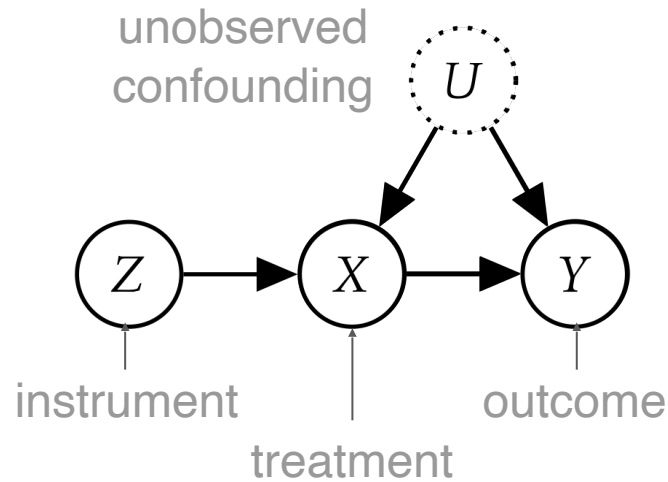


~~$$\mathbb{E}[e_Y | X] = 0$$~~

Losing hope...



Instrumental variables



(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 $\mathbb{E}[e_Y] = 0$

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

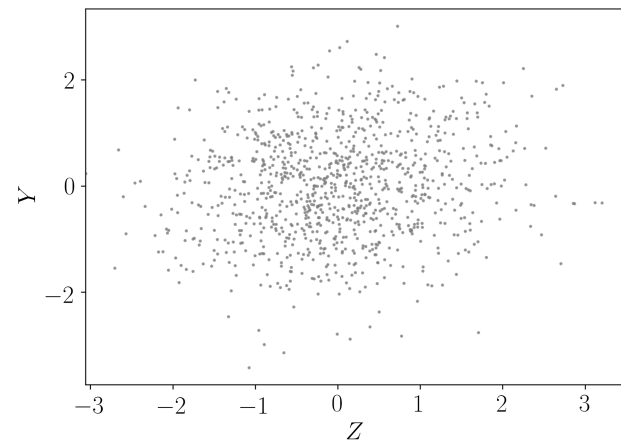
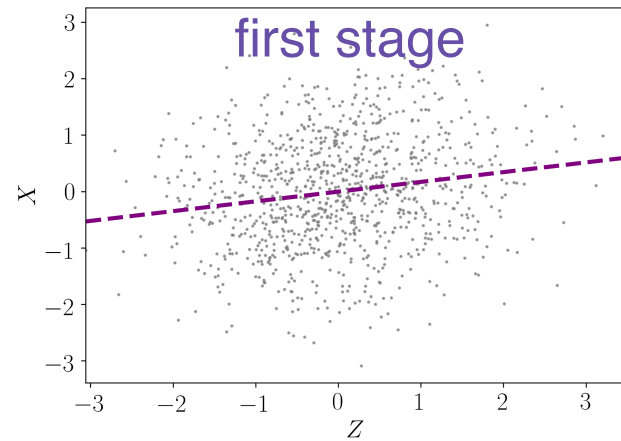
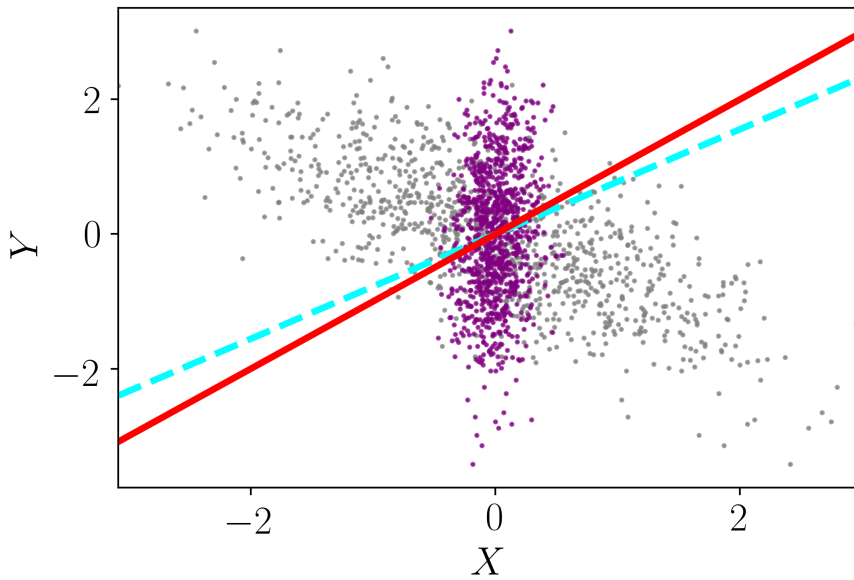
identifiable

unique under
mild conditions

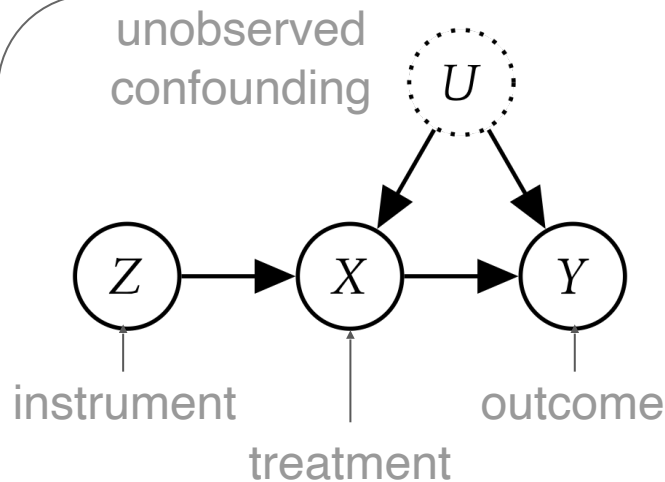
identifiable

Two stage least squares (2SLS) -- linear case

second stage



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 $Z \perp\!\!\!\perp Y \mid \{X, U\}$

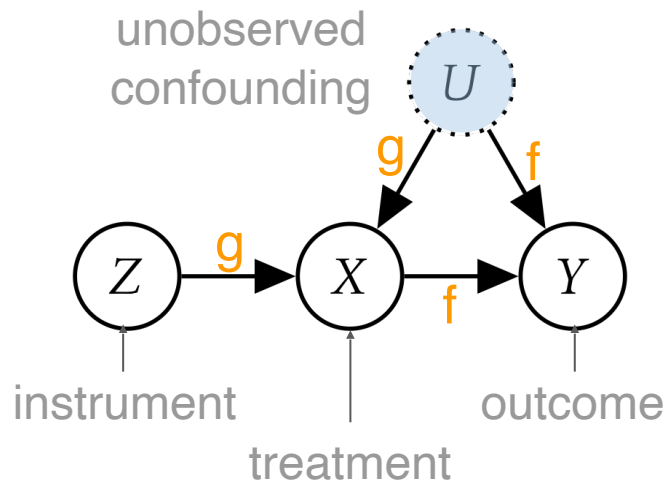
$$X = g(Z, U) \qquad Y = f(X, U)$$

non-linear, non-additive

Goal - partial identification

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

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



Equations defining the relationships:

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

Annotations:

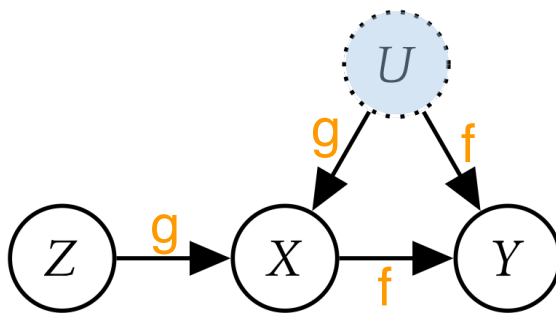
- Blue arrows point to U in both equations, labeled "optimize over 'all' distributions".
- An orange arrow points to g in the first equation, labeled "optimize over 'all' functions".
- An orange arrow points to f in the second equation, labeled "optimize over 'all' functions".

Goal

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

- 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**

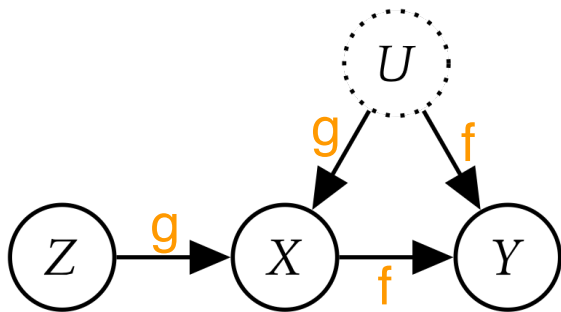
find convenient
representation of U from
which we can sample



choose convenient
function spaces

approximate constraints of
preserving $p(x | z)$ and $p(y | z)$

Our practical approach



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

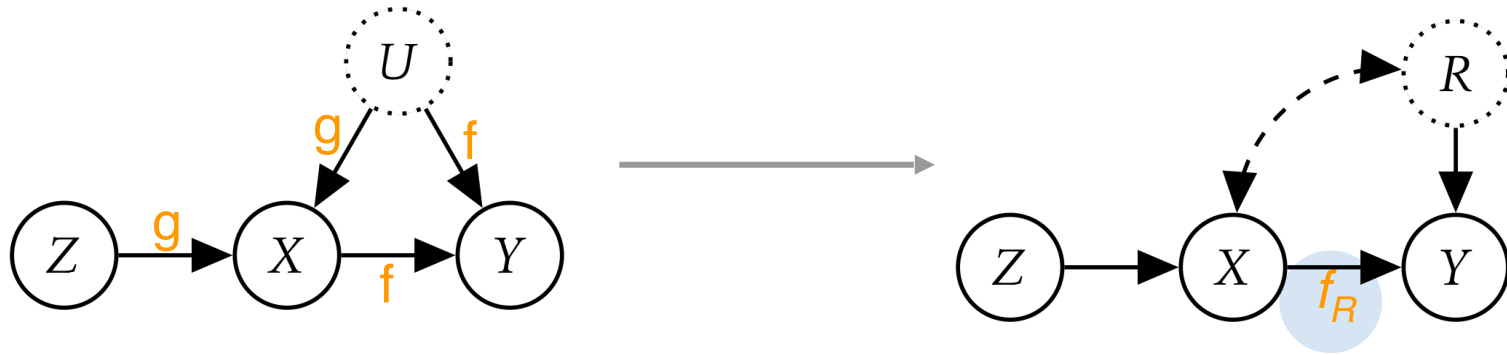
$$Y = f(X, U) = \lambda_1 X + \lambda_2 X U_1 + U_2$$

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

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

→ 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

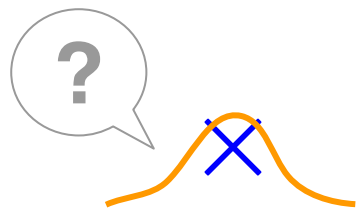


We choose a simple
parameterization

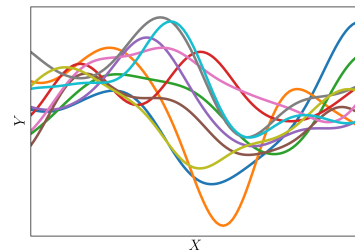
$$f_r(x) := f_{\theta_r}(x) \quad \text{for } \theta \in \Theta \subset \mathbb{R}^K$$

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

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



θ

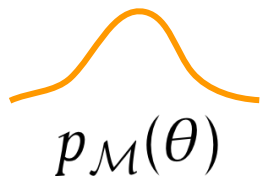


$f_{\theta} : X \rightarrow Y$

polynomials

neural networks

Gaussian process
samples



implies a causal model

Goal

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

$\int p_{\mathcal{M}}(x, y | z, \theta) p_{\mathcal{M}}(\theta) d\theta$ matches (estimated) marginals $p(x|z), p(y|z)$

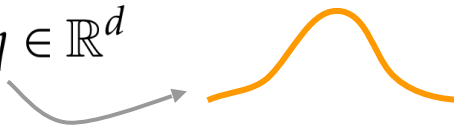
ideally

low variance Monte-Carlo
gradient estimation

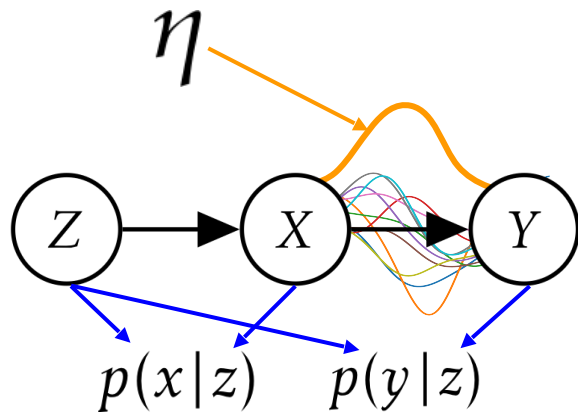
differentiable sampling

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

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



Objective function

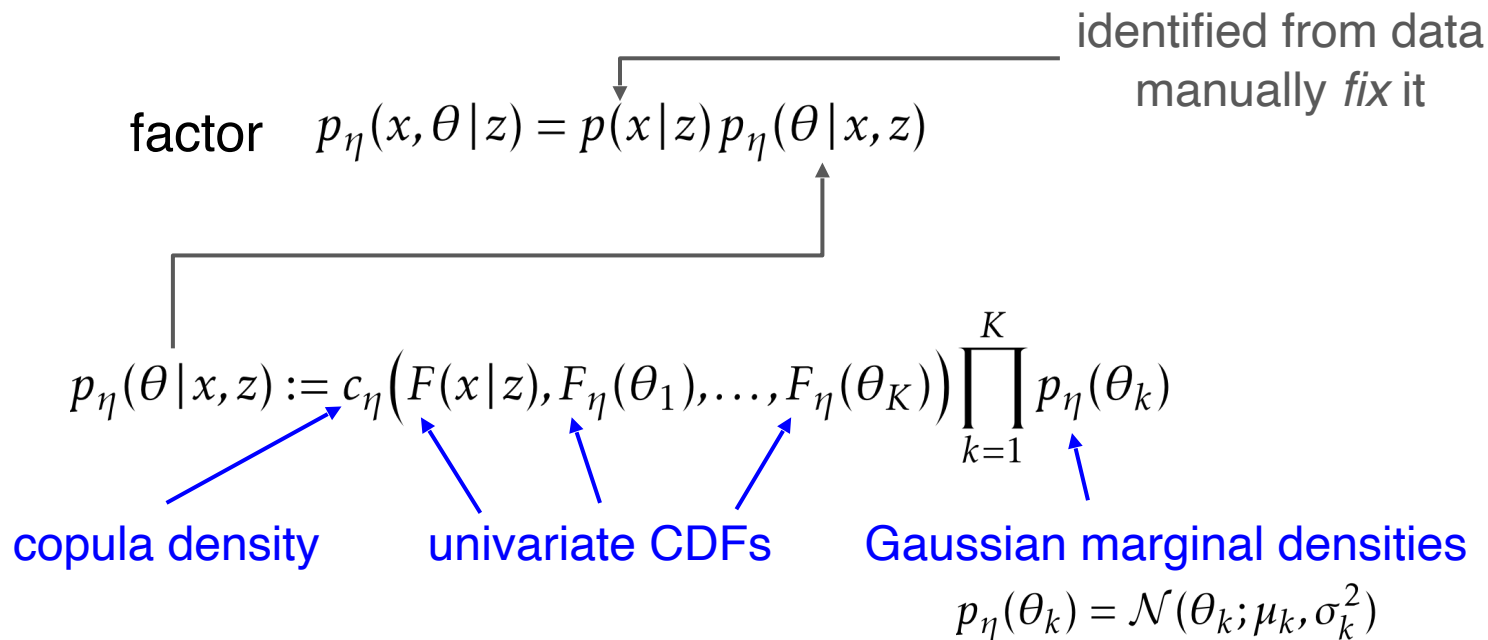


objective

$$\min_{\eta} / \max_{\eta} \mathbb{E}[Y | do(x^{\star})] = \min_{\eta} / \max_{\eta} \int f_{\theta}(x^{\star}) p_{\eta}(\theta) d\theta$$

Our model must match the observed data. Next up: Add these constraints.

Match $p(x|z)$ and enforcing $Z \perp U$



for a multivariate Gaussian copula, the optimization parameters are

$$\eta := \{\mu_1, \ln(\sigma_1^2), \dots, \mu_K, \ln(\sigma_K^2), L\} \in \mathbb{R}^{K(K+1)/2+2K}$$

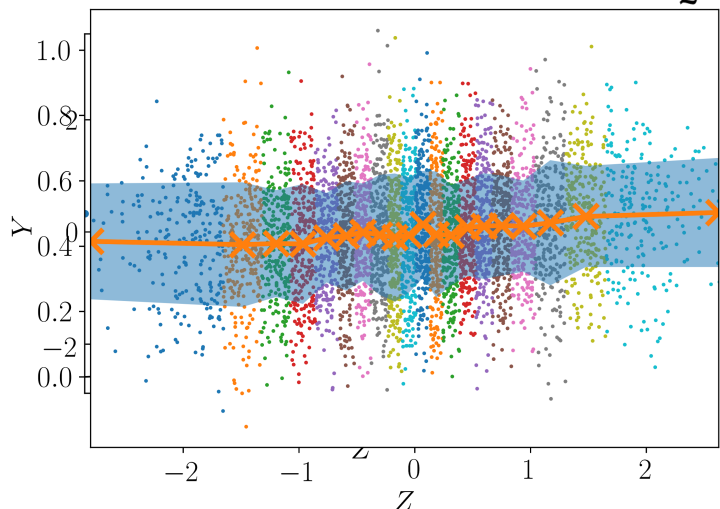
Match $p(y | z)$

exact constraint in the continuous outcome setting

data $\Pr(Y \leq y | Z = z) = \int \mathbf{1}(f_{\theta}(x) \leq y) p_{\eta}(x, \theta | z) dx d\theta$ our model

choose discrete finite grid of assignment points to bins

- finite number of constraints
 - integral over non-continuous indicator
- $$z^{(m)} := F_Z^{-1}\left(\frac{m}{M+1}\right) \text{ for } m \in [M]$$



for a dictionary of basis functions $\{\phi_l\}_{l=1}^L$

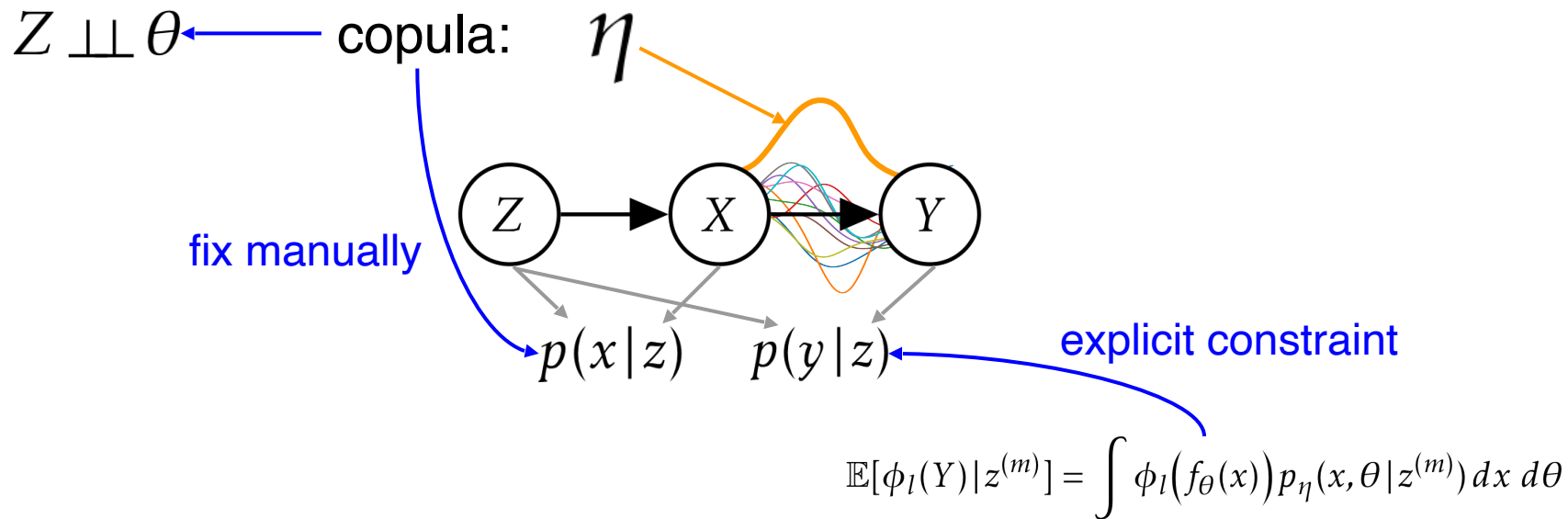
$$\mathbb{E}[\phi_l(Y) | z^{(m)}] = \int \phi_l(f_{\theta}(x)) p_{\eta}(x, \theta | z^{(m)}) dx d\theta$$

data

our model

$$\phi_1(Y) := \mathbb{E}[Y], \phi_2(Y) := \mathbb{V}[Y]$$

Intermediate overview



objective

$$\min_{\eta} / \max_{\eta} \mathbb{E}[Y | do(x^\star)] = \min_{\eta} / \max_{\eta} \int f_\theta(x^\star) p_\eta(\theta) d\theta$$

The final optimization problem

can sample from these in a differentiable fashion (w.r.t. η)

objective:

$$o_{x^\star}(\eta) := \int f_\theta(x^\star) p_\eta(\theta) d\theta$$

constraint LHS:

$$\text{LHS}_{m,l} := \mathbb{E}[\phi_l(Y) | z^{(m)}]$$

precompute once up front from data

constraint RHS:

$$\text{RHS}_{m,l}(\eta) := \int \phi_l(f_\theta(x)) p_\eta(x, \theta | z^{(m)}) dx d\theta$$

opt. problem:

$$\min_{\eta} / \max_{\eta} o_{x^\star}(\eta) \quad \text{s.t.} \quad \text{LHS}_{m,l} = \text{RHS}_{m,l}(\eta) \text{ for all } m \in [M], l \in [L]$$

only satisfy this approximately

use augmented Lagrangian with stochastic gradient descent

- for each $z^{(m)}$ sample batch of θ
- take average to estimate objective and constraint term RHS
- use auto-differentiation and gradient-based optimization



Empirical results

Choices of response functions

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

Polynomials

$$\psi_k(x) = x^{k-1} \text{ for } k \in [K]$$

Neural network

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

Gaussian process

Train GPs on subsets of observed data $X \rightarrow Y$ and take random samples from the GP as basis functions.

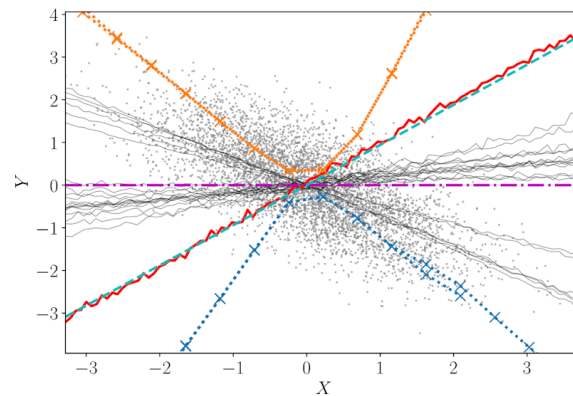
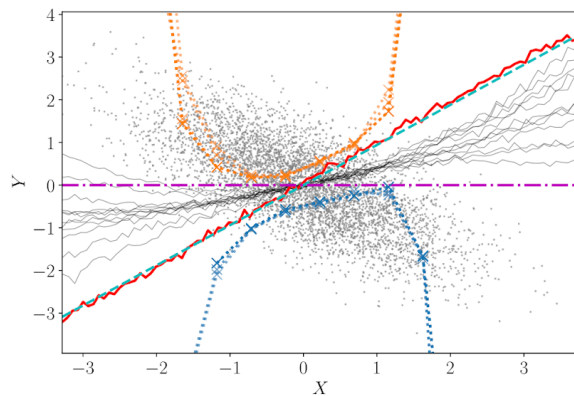
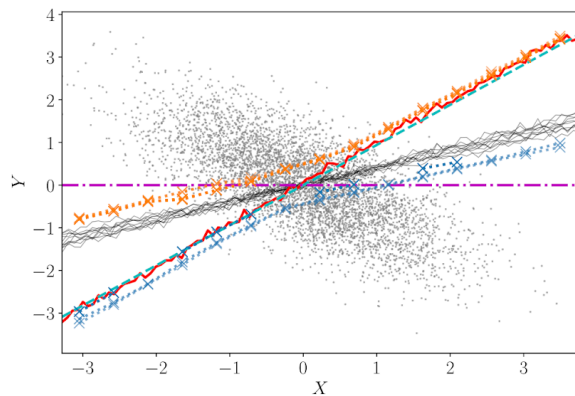
— possible models $E[Y|do(X = x^*)]$ - - - 2SLS - · - KIV · × · lower bound · × · upper bound · · · data

linear response

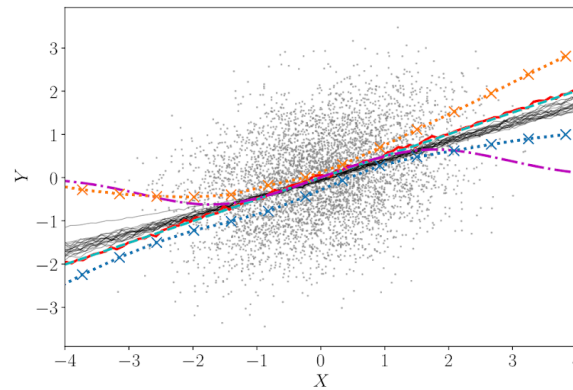
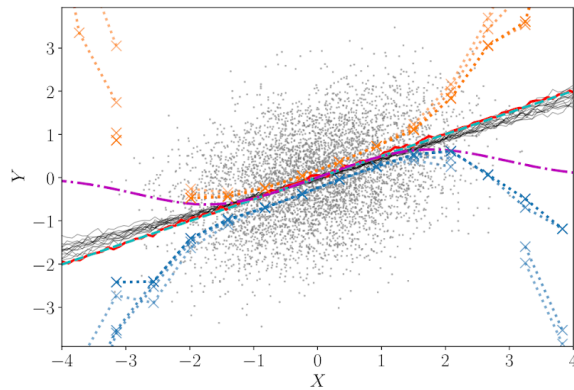
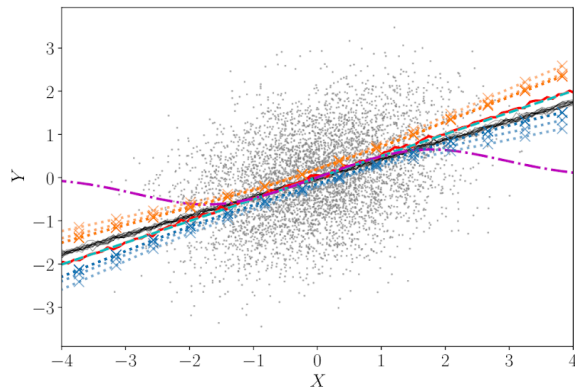
quadratic response

MLP response

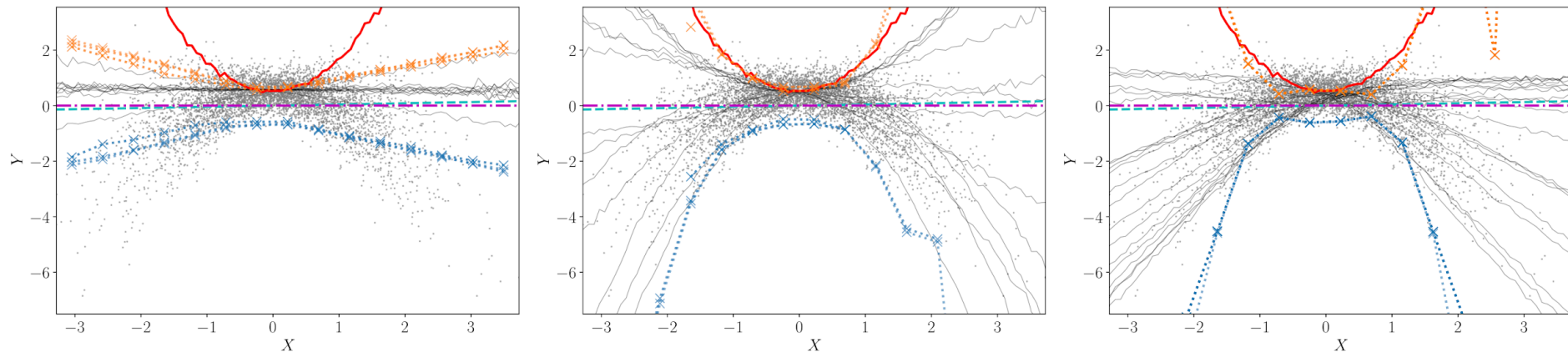
linear Gaussian setting; weak instrument and strong confounding ($\alpha = 0.5, \beta = 3$)



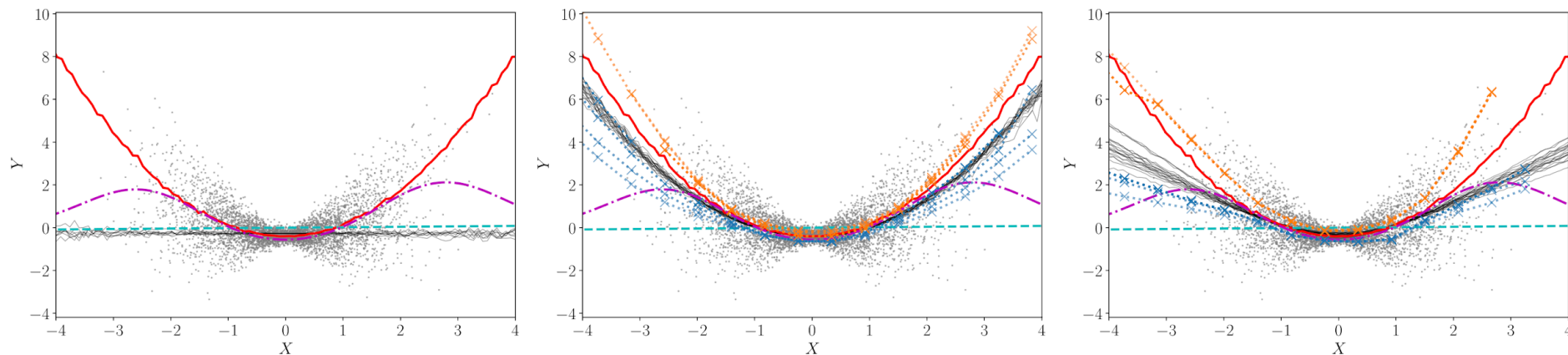
linear Gaussian setting; strong instrument and weak confounding ($\alpha = 3, \beta = 0.5$)



non-additive, non-linear setting; weak instrument and strong confounding ($\alpha = 0.5, \beta = 3$)

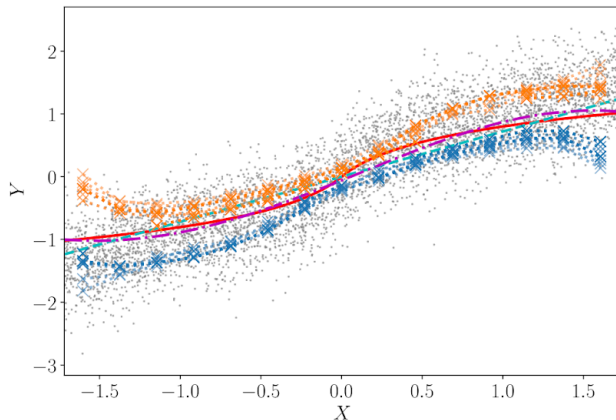


non-additive, non-linear setting; strong instrument and weak confounding ($\alpha = 3, \beta = 0.5$)

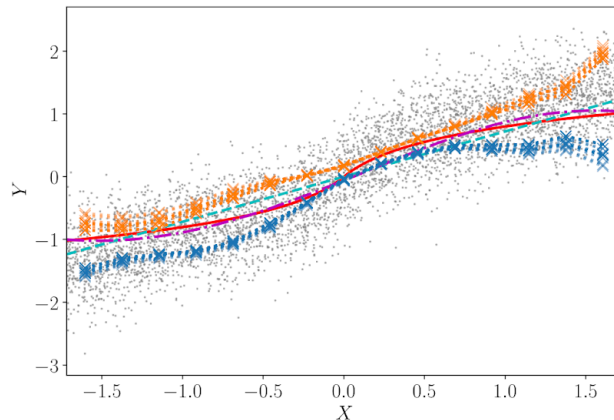


Sigmoidal cause-effect design

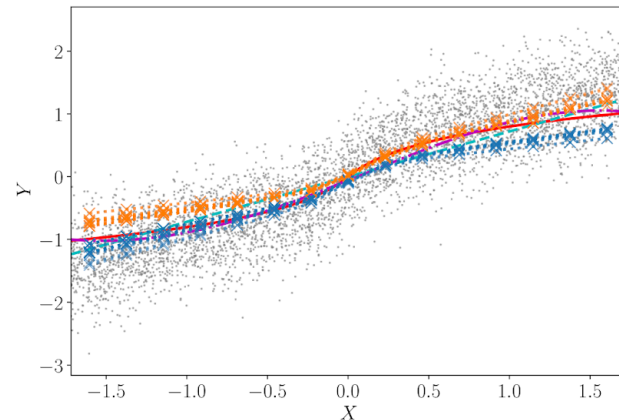
cubic response



GP response



MLP response



more details and experiments (also in the small data regime) in the paper

<https://arxiv.org/abs/2006.06366>

Thank you