Learning to act in noisy contexts using deep proxy learning

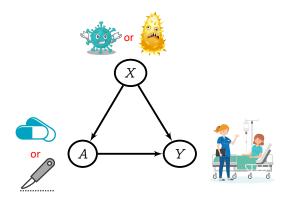
Arthur Gretton

Gatsby Computational Neuroscience Unit
Google Deepmind

Jump Trading CSML Series, 2025

Observation vs intervention

Conditioning from observation: $\mathbb{E}[Y|A=a] = \sum_{x} \mathbb{E}[Y|a,x] p(x|a)$



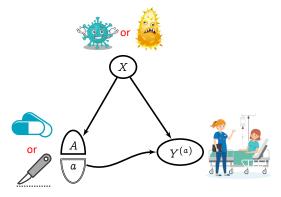
From our *observations* of historical hospital data:

- P(Y = cured | A = pills) = 0.85
- P(Y = cured|A = surgery) = 0.72

Observation vs intervention

Average causal effect/dose response curve (intervention):

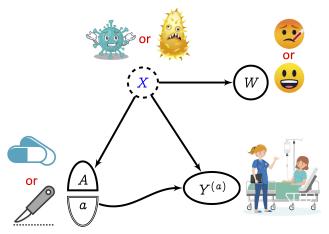
$$\mathbb{E}[Y^{(a)}] = \sum_x \mathbb{E}[Y|a,x] p(x)$$



From our *intervention* (making all patients take a treatment):

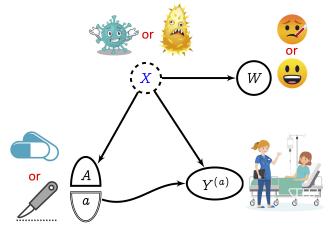
- $P(Y^{\text{(pills)}} = \text{cured}) = 0.64$
- $P(Y^{\text{(surgery)}} = \text{cured}) = 0.75$

We record symptom W, not disease X



- P(W = fever|X = mild) = 0.2
- P(W = fever|X = severe) = 0.8

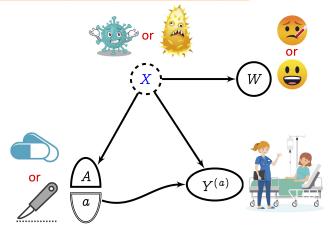
We record symptom W, not disease X



- P(W = fever | X = mild) = 0.2
- P(W = fever|X = severe) = 0.8

Could we just write: $P(Y^{(a)}) \stackrel{?}{=} \sum_{w \in \{0,1\}} \mathbb{E}[Y|a,w] p(w)$

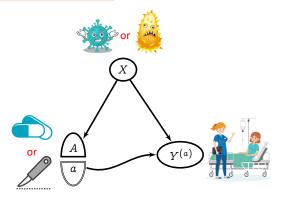
We record symptom W, not disease X



Wrong recommendation made:

Correct answer impossible without observing X

Some core assumptions



Assume:

- Stable Unit Treatment Value Assumption (aka "no interference"),
- Conditional exchangeability $Y^{(a)} \perp \!\!\!\perp A|X$.
- Overlap.

Outline

Causal effect estimation, with hidden covariates X:

■ Use proxy variables (negative controls)

Applications: effect of actions under

- privacy constraints (email, ads, DMA)
- data gathering constraints (edge computing)
- fundamental limitations (preferences, state of mind)

Outline

Causal effect estimation, with hidden covariates X:

■ Use proxy variables (negative controls)

Applications: effect of actions under

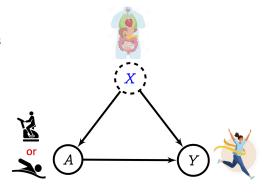
- privacy constraints (email, ads, DMA)
- data gathering constraints (edge computing)
- fundamental limitations (preferences, state of mind)

What's new and why?

- Treatment A, proxy variables, etc can be multivariate, complicated...
- ...by using adaptive neural net feature representations
- Don't meet your heroes model your hidden variables!

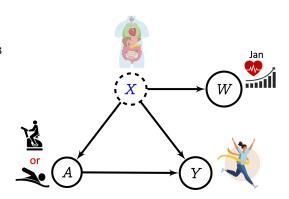
Unobserved X with (possibly) complex nonlinear effects on A, Y

- X: true physical status
- A: exercise regimes
- Y: fitness goal



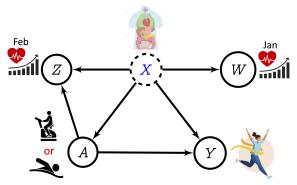
Unobserved X with (possibly) complex nonlinear effects on A, Y

- X: true physical status
- A: exercise regimes
- Y: fitness goal
- W: health readings before A



Unobserved X with (possibly) complex nonlinear effects on A, Y

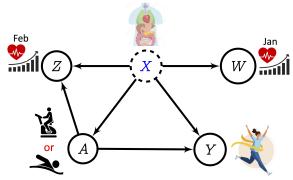
- X: true physical status
- *A*: exercise regimes
- Y: fitness goal
- W: health readings before A
- Z: health readings after A



Unobserved X with (possibly) complex nonlinear effects on A, Y

In this example:

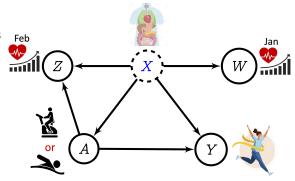
- X: true physical status
- A: exercise regimes
- Y: fitness goal
- W: health readings before A
- Z: health readings after A



 \implies Can recover $\mathbb{E}(Y^{(a)})$ from observational data

Unobserved X with (possibly) complex nonlinear effects on A, Y

- X: true physical status
- A: exercise regimes
- Y: fitness goal
- W: health readings before A
- Z: health readings after A

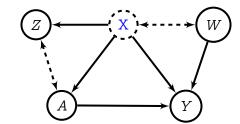


- \implies Can recover $\mathbb{E}(Y^{(a)})$ from observational data
- ⇒ More usefully: evaluate novel policy.

Proxy variables: general setting

Unobserved X with (possibly) complex nonlinear effects on A, Y. The definitions are:

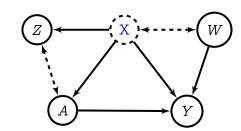
- X: unobserved confounder.
- A: treatment
- *Y*: outcome
- \blacksquare Z: treatment proxy
- W outcome proxy



Proxy variables: general setting

Unobserved X with (possibly) complex nonlinear effects on A, Y. The definitions are:

- X: unobserved confounder.
- *A*: treatment
- *Y*: outcome
- \blacksquare Z: treatment proxy
- W outcome proxy



Structural assumptions:

$$W \perp \!\!\!\perp (Z, A)|X$$

 $Y \perp \!\!\!\perp Z|(A, X)$

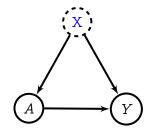
Miao, Geng, Tchetgen Tchetgen (2018): Identifying causal effects with proxy variables of an unmeasured confounder.

7/30

Why proxy variables? A simple proof

The definitions are:

- X: unobserved confounder.
- *A*: treatment
- *Y*: outcome



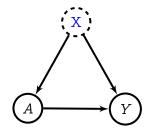
If X were observed,

$$\underbrace{P(Y^{(a)})}_{d_u imes 1} := \sum_{i=1}^{d_x} P(Y|oldsymbol{x}_i, \, a) P(oldsymbol{x}_i)$$

Why proxy variables? A simple proof

The definitions are:

- X: unobserved confounder.
- A: treatment
- *Y*: outcome



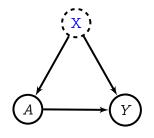
If X were observed,

$$\underbrace{P(Y^{(a)})}_{d_y imes 1} := \sum_{i=1}^{d_x} P(Y|x_i, a) P(x_i) = \underbrace{P(Y|X, a)}_{d_y imes d_x} \underbrace{P(X)}_{d_x imes 1}$$

Why proxy variables? A simple proof

The definitions are:

- X: unobserved confounder.
- A: treatment
- *Y*: outcome



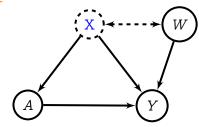
If X were observed,

$$\underbrace{P(\,Y^{(a)})}_{d_y imes 1} := \sum_{i=1}^{d_x} P(\,Y|oldsymbol{x}_i,\,a) P(oldsymbol{x}_i) = \underbrace{P(\,Y|oldsymbol{X},\,a) P(oldsymbol{X})}_{d_y imes d_x} \underbrace{P(\,Y|oldsymbol{X},\,a) P(oldsymbol{X})}_{d_x imes 1}$$

Goal: "get rid of the blue" X

The definitions are:

- X: unobserved confounder.
- *A*: treatment
- *Y*: outcome
- W: outcome proxy

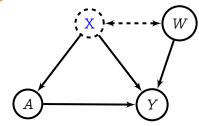


For each a, if we could solve:

$$\underbrace{P(\mathit{Y}|\mathit{X},\mathit{a})}_{\mathit{d_{\mathit{y}}} imes\mathit{d_{\mathit{x}}}} = \underbrace{\mathit{H}_{w,\mathit{a}}}_{\mathit{d_{\mathit{y}}} imes\mathit{d_{\mathit{w}}}}\underbrace{P(\mathit{W}|\mathit{X})}_{\mathit{d_{\mathit{w}}} imes\mathit{d_{\mathit{x}}}}$$

The definitions are:

- *X*: unobserved confounder.
- A: treatment
- *Y*: outcome
- W: outcome proxy



For each a, if we could solve:

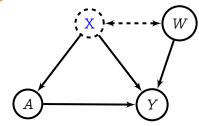
$$\underbrace{P(Y|X,a)}_{d_y imes d_x} = \underbrace{H_{w,a}}_{d_y imes d_w} \underbrace{P(W|X)}_{d_w imes d_x}$$

.....then

$$P(Y^{(a)}) = P(Y|X, a)P(X)$$

The definitions are:

- X: unobserved confounder.
- A: treatment
- *Y*: outcome
- W: outcome proxy



For each a, if we could solve:

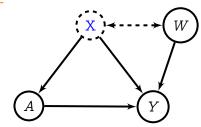
$$\underbrace{P(Y|X,a)}_{d_y imes d_x} = \underbrace{H_{w,a}}_{d_y imes d_w} \underbrace{P(W|X)}_{d_w imes d_x}$$

.....then

$$P(Y^{(a)}) = P(Y|X, a)P(X)$$
$$= H_{w,a}P(W|X)P(X)$$

The definitions are:

- *X*: unobserved confounder.
- A: treatment
- *Y*: outcome
- W: outcome proxy



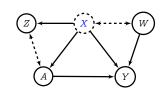
For each a, if we could solve:

$$\underbrace{P(\mathit{Y}|\mathit{X},\mathit{a})}_{\mathit{d_{\mathit{y}}} imes\mathit{d_{\mathit{x}}}} = \underbrace{\mathit{H_{w,a}}}_{\mathit{d_{\mathit{y}}} imes\mathit{d_{\mathit{w}}}} \underbrace{P(\mathit{W}|\mathit{X})}_{\mathit{d_{\mathit{w}}} imes\mathit{d_{\mathit{x}}}}$$

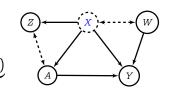
.....then

$$egin{aligned} P(Y^{(a)}) &= P(Y|X,a)P(X) \ &= H_{w,a}P(W|X)P(X) \ &= H_{w,a}P(W) \end{aligned}$$

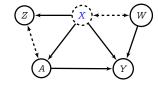
$$P(Y|X,a) = H_{w,a}P(W|X)$$



$$P(Y|X,a) \underbrace{p(X|Z,a)}_{d_x imes d_z} = H_{w,a} P(W|X) \underbrace{p(X|Z,a)}_{d_x imes d_z}$$



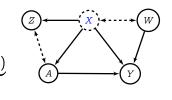
$$P(Y|X,a) \underbrace{p(X|Z,a)}_{d_x imes d_z} = H_{w,a} P(W|X) \underbrace{p(X|Z,a)}_{d_x imes d_z}$$



Because
$$W \perp\!\!\!\perp (Z,A)|X,$$

$$P(W|X)p(X|Z,a) = P(W|Z,a)$$

$$P(Y|X,a)\underbrace{p(X|Z,a)}_{d_x \times d_z} = H_{w,a}P(W|X)\underbrace{p(X|Z,a)}_{d_x \times d_z}$$



Because
$$W \perp \!\!\!\perp (Z, A)|X$$
,

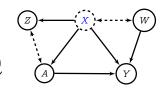
$$P(W|X)p(X|Z,a) = P(W|Z,a)$$

Because
$$Y \perp \!\!\!\perp Z | (A, X)$$
,

$$P(Y|X,a)p(X|Z,a) = P(Y|Z,a)$$

From last slide,

$$P(Y|X,a) \underbrace{p(X|Z,a)}_{d_x \times d_z} = H_{w,a} P(W|X) \underbrace{p(X|Z,a)}_{d_x \times d_z}$$



Because
$$W \perp \!\!\!\perp (Z, A)|X$$
,

$$P(W|X)p(X|Z,a) = P(W|Z,a)$$

Because
$$Y \perp \!\!\!\perp Z | (A, X)$$
,

$$P(Y|X,a)p(X|Z,a) = P(Y|Z,a)$$

Solve for $H_{w,a}$:

$$P(Y|Z,a) = H_{w,a}P(W|Z,a)$$

Everything observed!

Proxy/Negative Control Methods in the Real World

Unobserved confounders: proxy methods

Kernel features (ICML 2021):









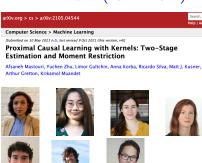




Code for NN and kernel proxy methods:

Unobserved confounders: proxy methods

Kernel features (ICML 2021):





Code for NN and kernel proxy methods:

https://github.com/liyuan9988/DeepFeatureProxyVariable/

One model: linear functions of features

All learned functions will take the form:

$$oldsymbol{\gamma}(x) = oldsymbol{\gamma}^ op arphi_{ heta}(x)$$

NN approach: Finite dictionaries of learned neural net features $\varphi_{\theta}(x)$ (linear final layer γ)

Xu, G., A Neural mean embedding approach for back-door and front-door adjustment. (ICLR 23) Xu, Chen, Srinivasan, de Freitas, Doucet, G. Learning Deep Features in Instrumental Variable Regression. (ICLR 21)

Xu, Kanagawa, G. "Deep Proxy Causal Learning and its Application to Confounded Bandit Policy Evaluation". (NeurIPS 21)

Model fitting: neural ridge regression

Learn $\gamma_0(x) := \mathbb{E}[\,Y|X=x]$ from features $arphi_{ heta}(x_i)$ with outcomes y_i :

$$\hat{\gamma} = \arg\min_{\gamma \in \mathcal{H}} \left(\sum_{i=1}^{n} \left(y_i - \gamma^{\top} \varphi_{\theta}(x_i) \right)^2 + \lambda \|\gamma\|^2 \right)$$
 (1)

Model fitting: neural ridge regression

Learn $\gamma_0(x) := \mathbb{E}[Y|X=x]$ from features $\varphi_{\theta}(x_i)$ with outcomes y_i :

$$\hat{\gamma} = \arg\min_{\gamma \in \mathcal{H}} \left(\sum_{i=1}^{n} \left(y_i - \gamma^{\top} \varphi_{\theta}(x_i) \right)^2 + \lambda \|\gamma\|^2 \right)$$
 (1)

Solution for linear final layer γ :

$$egin{aligned} \hat{\gamma} &= C_{YX}^{(heta)} (\, C_{XX}^{(heta)} + \lambda)^{-1} \ C_{YX}^{(heta)} &= rac{1}{n} \sum_{i=1}^n [y_i \; arphi_{ heta}(x_i)^ op] \ C_{XX}^{(heta)} &= rac{1}{n} \sum_{i=1}^n [arphi_{ heta}(x_i) \; arphi_{ heta}(x_i)^ op] \end{aligned}$$

Model fitting: neural ridge regression

Learn $\gamma_0(x) := \mathbb{E}[Y|X=x]$ from features $\varphi_{\theta}(x_i)$ with outcomes y_i :

$$\hat{\gamma} = \arg\min_{\gamma \in \mathcal{H}} \left(\sum_{i=1}^n \left(y_i - \gamma^\top \varphi_{\theta}(x_i) \right)^2 + \lambda \|\gamma\|^2 \right)$$
 (1)

Solution for linear final layer γ :

$$egin{aligned} \hat{\gamma} &= C_{YX}^{(heta)} (\, C_{XX}^{(heta)} + \lambda)^{-1} \ C_{YX}^{(heta)} &= rac{1}{n} \sum_{i=1}^n [y_i \ arphi_{ heta}(x_i)^ op] \ C_{XX}^{(heta)} &= rac{1}{n} \sum_{i=1}^n [arphi_{ heta}(x_i) \ arphi_{ heta}(x_i) \ arphi_{ heta}(x_i)^ op] \end{aligned}$$

How to solve for θ :

Substitute $\hat{\gamma}$ into (1), backprop through Cholesky for θ .

More details: Galashov, Da Costa, Xu, Hennig, G, Closed-Form Last Layer Optimization (2025, arxiv:2510.04606)

Road map: NN proxy learning

We'll proceed as follows:

- Proxy relation for continuous variables
- Loss function for deep proxy learning
- Define primary (ridge) regression with this loss
- Define secondary (ridge) regression as input to primary

If X were observed, we would write (dose-response curve)

$$\mathbb{E}(Y^{(a)}) = \int_x \mathbb{E}(Y|a,x)p(x)dx.$$

....but we do not observe X.

If X were observed, we would write (dose-response curve)

$$\mathbb{E}(Y^{(a)}) = \int_x \mathbb{E}(Y|a,x)p(x)dx.$$

....but we do not observe X.

Main theorem: Assume we solved for link function:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

- "Primary" $\mathbb{E}(Y|a,z)$, "secondary" $\mathbb{E}_{W|a,z}$ linked by h_y
- All variables observed, X not seen or modeled.

Fredholm equation of first kind. Link existence requires \Diamond , identification of ATE requires \triangle (and further technical assumptions) [XKG: Asspumption 2, Prop. 1,Corr. 1; Deaner]

$$\mathbb{E}[f(X)|A=a,Z=z]=0, \ \forall (z,a) \iff f(X)=0, \ \mathbb{P}_X \text{ a.s.} \quad \triangle$$

$$\mathbb{E}[f(X)|A=a,W=w]=0, \ \forall (w,a) \iff f(X)=0, \ \mathbb{P}_X \text{ a.s.} \quad \diamondsuit$$
17/30

If X were observed, we would write (dose-response curve)

$$\mathbb{E}(Y^{(a)}) = \int_x \mathbb{E}(Y|a,x)p(x)dx.$$

....but we do not observe X.

Main theorem: Assume we solved for link function:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

- "Primary" $\mathbb{E}(Y|a,z)$, "secondary" $\mathbb{E}_{W|a,z}$ linked by h_y
- \blacksquare All variables observed, X not seen or modeled.

Dose-response curve via p(w):

$$\mathbb{E}(\mathit{Y}^{(a)}) = \int_{w} \mathit{h}_{y}(a,w) p(w) dw$$

If X were observed, we would write (dose-response curve)

$$\mathbb{E}(Y^{(a)}) = \int_x \mathbb{E}(Y|a,x)p(x)dx.$$

....but we do not observe X.

Main theorem: Assume we solved for link function:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

- "Primary" $\mathbb{E}(Y|a,z)$, "secondary" $\mathbb{E}_{W|a,z}$ linked by h_y
- \blacksquare All variables observed, X not seen or modeled.

Dose-response curve via p(w):

$$\mathbb{E}(\mathit{Y}^{(a)}) = \int_{w} \mathit{h}_{y}(a,w) p(w) dw$$

Challenge: need a loss function for h_y

Primary loss function for $h_y(w, a)$

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Primary loss function:

$$\hat{h}_{y} = rg \min_{h_{y}} \mathbb{E}_{Y,A,Z} \left(Y - \mathbb{E}_{W|A,Z} h_{y}(W,A) \right)^{2}$$

Why?

Primary loss function for $h_y(w, a)$

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z} h_y(W,a)$$

Primary loss function:

$$\hat{h}_{y} = rg\min_{h_{y}} \mathbb{E}_{Y,A,Z} \left(Y - \mathbb{E}_{W|A,Z} h_{y}(W,A) \right)^{2}$$

Why?

$$f^*(a,z) = \mathbb{E}(\,Y|\,a,z) ext{ solves}
onumber \ rgmin_f \mathbb{E}_{\,Y,A,Z} \,(\,Y-f(A,Z))^2$$

```
Deaner (2021).
Mastouri, Zhu, Gultchin, Korba, Silva, Kusner, G., Muandet (2021).
Xu, Kanagawa, G. (2021).
```

Primary loss function for $h_y(w, a)$

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z} h_y(W,a)$$

Primary loss function:

$$\hat{h}_{y} = rg \min_{h_{y}} \mathbb{E}_{Y,A,Z} \left(Y - \mathbb{E}_{W|A,Z} h_{y}(W,A) \right)^{2}$$

Why?

$$f^*(a,z) = \mathbb{E}(\,Y|a,z) ext{ solves}
onumber \ rgmin_f \mathbb{E}_{\,Y,A,Z} \, (\,Y-f(A,Z))^2$$

...and by the proxy model above,

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Deaner (2021). Mastouri, Zhu, Gultchin, Korba, Silva, Kusner, G., Muandet (2021). Xu, Kanagawa, G. (2021).

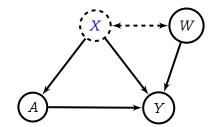
NN for link $h_y(a, w)$

The link function is a function of two arguments

$$h_y(a,w) = \gamma^ op \left[arphi_{ heta}(w) \otimes arphi_{\xi}(a)
ight] = \gamma^ op \left[egin{array}{c} arphi_{ heta,1}(w) arphi_{\xi,1}(a) \ arphi_{ heta,1}(w) arphi_{\xi,2}(a) \ dots \ \ dots \ dots \ dots \ dots \ dots \ dots \ \ dots \ dots \$$

Assume we have:

- output proxy NN features $\varphi_{\theta}(w)$
- lacksquare treatment NN features $arphi_{\xi}(a)$
- linear final layer γ
 (argument of feature map indicates feature space)



NN for link $h_y(a, w)$

The link function is a function of two arguments

$$h_y(a,w) = \gamma^ op \left[arphi_ heta(w) \otimes arphi_\xi(a)
ight]$$

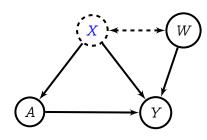
Assume we have:

- lacksquare output proxy NN features $\varphi_{\theta}(w)$
- lacktriangle treatment NN features $arphi_{\xi}(a)$
- In linear final layer γ (argument of feature map indicates feature space)

Questions:

- Why feature map $\varphi_{\theta}(w) \otimes \varphi_{\xi}(a)$?
- Why final linear layer γ ?

Both are necessary (next slide)!



Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Primary regression:

$$\hat{h}_y = rg\min_{h_y} \mathbb{E}_{Y,A,Z} \left(Y - \mathbb{E}_{oldsymbol{W}|A,Z} h_y(oldsymbol{W},A)
ight)^2 + \lambda_2 \|\gamma\|^2$$

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Primary regression:

$$\hat{h}_y = \arg\min_{h_y} \mathbb{E}_{Y,A,Z} \left(\left. Y - \mathbb{E}_{\left. W \right| A,Z} h_y (\left. \begin{matrix} W \end{matrix}, A \right) \right)^2 + \lambda_2 \| \gamma \|^2$$

How to get conditional expectation $\mathbb{E}_{W|a,z}h_y(W,a)$?

Density estimation for p(W|a, z)? Sample from p(W|a, z)?

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Primary regression:

$$\hat{h}_y = \arg\min_{h_y} \mathbb{E}_{Y,A,Z} \left(\left. Y - \mathbb{E}_{\boldsymbol{W}|A,Z} h_y(\boldsymbol{W},A) \right)^2 + \lambda_2 \|\gamma\|^2 \right.$$

Recall link function

$$h_y(extit{W}, a) = \left[\gamma^ op (arphi_ heta(extit{W}) \otimes arphi_\xi(a))
ight]$$

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Primary regression:

$$\hat{h}_y = \arg\min_{h_y} \mathbb{E}_{Y,A,Z} \left(\left. Y - \mathbb{E}_{\boldsymbol{W}|A,Z} h_y(\boldsymbol{W},A) \right)^2 + \lambda_2 \|\gamma\|^2 \right.$$

Recall link function

$$\mathbb{E}_{W|a,z} \; h_y(\hspace{.05cm} W,\hspace{.05cm} a) = \hspace{.05cm} \mathbb{E}_{\hspace{.05cm} W|a,z} \; \left[\gamma^{ op} (\hspace{.05cm} arphi_{ heta}(\hspace{.05cm} W) \otimes arphi_{\xi}(\hspace{.05cm} a))
ight]$$

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Primary regression:

$$\hat{h}_y = \arg\min_{h_y} \mathbb{E}_{Y,A,Z} \left(\left. Y - \mathbb{E}_{\left. W \right| A,Z} h_y (\left. W \right, A) \right)^2 + \lambda_2 \| \gamma \|^2$$

Recall link function

$$egin{aligned} \mathbb{E}_{W|a,z} \; h_y(\,W,\,a) &= \; \mathbb{E}_{W|a,z} \; \left[\gamma^ op \left(arphi_ heta(\,W) \otimes arphi_ ext{\xi}(a)
ight)
ight] \ &= \gamma^ op \left(\mathbb{E}_{W|a,z} \left[arphi_ heta(\,W)
ight] \otimes arphi_ ext{\xi}(a)
ight) \end{aligned}$$

(this is why linear γ and feature map $\varphi_{\theta}(w) \otimes \varphi_{\xi}(a)$)

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Primary regression:

$$\hat{h}_y = \arg\min_{h_y} \mathbb{E}_{Y,A,Z} \left(\left. Y - \mathbb{E}_{\left. W \right| A,Z} h_y (\left. W \right, A) \right)^2 + \lambda_2 \| \gamma \|^2$$

Recall link function

$$egin{aligned} \mathbb{E}_{W|a,z} \; h_y(\,W,\,a) &= \; \mathbb{E}_{W|a,z} \; \left[\gamma^ op \left(arphi_ heta(\,W) \otimes arphi_\xi(a)
ight)
ight] \ &= \gamma^ op \left(\mathbb{E}_{W|a,z} \left[arphi_ heta(\,W)
ight] \otimes arphi_\xi(a)
ight) \ & ext{cond. feat. mean} \end{aligned}$$

Ridge regression (again!)

$$\mathbb{E}_{W|a,z}arphi_{ heta}(W)=\hat{F}_{ heta,\zeta}arphi_{\zeta}(a,z)$$

NN ridge regression for $\mathbb{E}_{W|a,z}\varphi_{\theta}(W)$

Secondary regression: learn NN features $\phi_{\zeta}(Z)$ and linear layer F:

$$\mathbb{E}_{W|a,z}arphi_{ heta}(W)=\hat{F}_{ heta,\zeta}arphi_{\zeta}(a,z)$$

with RR loss

$$\mathbb{E}_{W,A,Z} \left\| arphi_{ heta}(W) - rac{oldsymbol{F}}{oldsymbol{F}} arphi_{\zeta}(A,Z)
ight\|^2 + \lambda_1 \|rac{oldsymbol{F}}{oldsymbol{F}} \|^2$$

 $\hat{F}_{\theta,\zeta}$ in closed form wrt $\phi_{\theta},\phi_{\zeta}$.

NN ridge regression for $\mathbb{E}_{W|a,z}\varphi_{\theta}(W)$

Secondary regression: learn NN features $\phi_{\zeta}(Z)$ and linear layer F:

$$\mathbb{E}_{W|a,z}arphi_{ heta}(W)=\hat{F}_{ heta,\zeta}arphi_{\zeta}(a,z)$$

with RR loss

$$\mathbb{E}_{W,A,Z} \left\| arphi_{ heta}(W) - rac{F}{F} arphi_{\zeta}(A,Z)
ight\|^2 + \lambda_1 \|rac{F}{F}\|^2$$

 $\hat{F}_{\theta,\zeta}$ in closed form wrt $\phi_{\theta},\phi_{\zeta}$.

Plug $\hat{F}_{\theta,\zeta}$ into S1 loss, backprop through Cholesky for ζ (...not θ ...why not?)

Solve for θ, ξ, ζ :

Repeat until convergence:

■ Secondary: Solve for $\hat{F}_{\theta,\zeta}$, then gradient steps on ζ (backprop through Cholesky)

```
Solve for \theta, \xi, \zeta:
```

Repeat until convergence:

- Secondary: Solve for $\hat{F}_{\theta,\zeta}$, then gradient steps on ζ (backprop through Cholesky)
- Primary: Solve for $\hat{\gamma}$ in terms of $\hat{F}_{\theta,\zeta}\phi_{\zeta}(A,Z)$ and $\varphi_{\xi}(A)$

```
Solve for \theta, \xi, \zeta:
```

Repeat until convergence:

- Secondary: Solve for $\hat{F}_{\theta,\zeta}$, then gradient steps on ζ (backprop through Cholesky)
- Primary: Solve for $\hat{\gamma}$ in terms of $\hat{F}_{\theta,\zeta}\phi_{\zeta}(A,Z)$ and $\varphi_{\xi}(A)$
- Primary: Gradient steps on θ, ξ (backprop through Cholesky)
 - $\hat{F}_{\theta,\zeta}$ remains optimal wrt current φ_{θ} .

```
Solve for \theta, \xi, \zeta:
```

Repeat until convergence:

- Secondary: Solve for $\hat{F}_{\theta,\zeta}$, then gradient steps on ζ (backprop through Cholesky)
- Primary: Solve for $\hat{\gamma}$ in terms of $\hat{F}_{\theta,\zeta}\phi_{\zeta}(A,Z)$ and $\varphi_{\xi}(A)$
- Primary: Gradient steps on θ, ξ (backprop through Cholesky)
 - $\hat{F}_{\theta,\zeta}$ remains optimal wrt current φ_{θ} .

Iterate between updates of θ , ξ and ζ

Solve for θ, ξ, ζ :

Repeat until convergence:

- Secondary: Solve for $\hat{F}_{\theta,\zeta}$, then gradient steps on ζ (backprop through Cholesky)
- Primary: Solve for $\hat{\gamma}$ in terms of $\hat{F}_{\theta,\zeta}\phi_{\zeta}(A,Z)$ and $\varphi_{\xi}(A)$
- Primary: Gradient steps on θ, ξ (backprop through Cholesky)
 - $\hat{F}_{\theta,\zeta}$ remains optimal wrt current φ_{θ} .

Iterate between updates of θ , ξ and ζ

Key point: features $\varphi_{\theta}(W)$ learned specially for:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

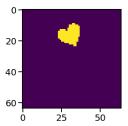
Contrast with autoencoders/sampling: must reconstruct/sample all of W.

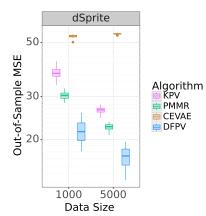
Experiments

Synthetic experiment, adaptive neural net features

dSprite example:

- $X = \{ scale, rotation, posX, posY \}$
- Treatment A is the image generated (with Gaussian noise)
- Outcome Y is quadratic function of A with multiplicative confounding by posY.
- Z = {scale, rotation, posX}, W = noisy image sharing posY
- Comparison with CEVAE (Louzios et al. 2017)



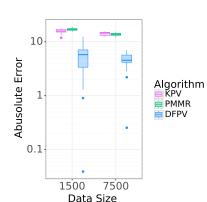


 ${\tt Louizos, Shalit, Mooij, Sontag, Zemel, Welling, Causal\ Effect\ Inference\ with\ Deep\ Latent-Variable} {\tt 24/30}\ {\tt Models\ (2017)}$

Confounded offline policy evaluation

Synthetic dataset, demand prediction for flight purchase.

- Treatment A is ticket price.
- Policy $A \sim \pi(Z)$ depends on fuel price.

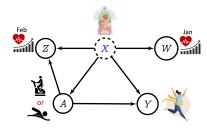


Conclusion

Causal effect estimation with unobserved X, (possibly) complex nonlinear effects on A, Y

We need to observe:

- Treatment proxy Z (interacts with A, but not directly with Y)
- Outcome proxy W (no direct interaction with A, can affect Y)



Conclusion

Causal effect estimation with unobserved X, (possibly) complex nonlinear effects on A, Y

We need to observe:

- Treatment proxy Z (interacts with A, but not directly with Y)
- Outcome proxy W (no direct interaction with A, can affect Y)

Feb Z X W Jan W Y Y

Key messages:

- Don't meet your heroes model/sample latents X
- \blacksquare Don't model all of W, only relevant features for Y
- "Ridge regression is all you need"

Code available:

https://github.com/liyuan9988/DeepFeatureProxyVariable/

Research support

Work supported by:

The Gatsby Charitable Foundation



Google Deepmind



Questions?



Failures of completeness assumptions (1)

Recall (one of the) completeness assumptions:

$$\mathbb{E}[f(X)|A=a,Z=z]=0, \forall (a,z) \iff f(X)=0, \mathbb{P}_X \text{ a.s.} \quad (\triangle)$$

For conciseness, assume conditioning on some a.

Failure 1: $Z \perp \!\!\! \perp X$ (no information about X in proxy)

$$egin{aligned} g(oldsymbol{X}|) &= ilde{g}(oldsymbol{X}) - \mathbb{E}_{oldsymbol{X}} ilde{g}(oldsymbol{X}) \ \mathbb{E}(g(oldsymbol{X})|oldsymbol{Z}, \, a) &= \mathbb{E}g(oldsymbol{X}) = 0. \end{aligned}$$

Failures of identifiability assumptions (2)

Failure 2: "exploitable invariance" of p(X|z)

$$egin{aligned} X &\sim \mathcal{N}(0,1), \ Z &= |X| + \mathcal{N}(0,1), \end{aligned}$$

where $p(X|z) \propto p(z|X)p(X)$ symmetric in X. Consider square integrable antisymmetric function $g(X) = -g(-X) \neq 0$. Then

$$egin{aligned} \mathbb{E}[g(X)|Z=z] &= \int_{-\infty}^{\infty} g(X)p(X|z)dX \ &= \int_{-\infty}^{0} g(X)p(X|z)dX + \int_{0}^{\infty} g(X)p(X|z)dX \ &= 0. \end{aligned}$$

If distribution of X|Z retains the same "symmetry class" over a set of Z with nonzero measure, then the assumption is violated by g(X) with zero mean on this class.